# Penetration Vision through Virtual Reality Headsets: Identifying 360-degree Videos from Head Movements

Anh Nguyen
*George Mason University*

Xiaokuan Zhang
*George Mason University*

Zhisheng Yan
*George Mason University*

arXiv:2402.11446v2 [cs.HC] 8 Mar 2024

## Abstract

In this paper, we present the *first* contactless side-channel attack for identifying 360° videos being viewed in a Virtual Reality (VR) Head Mounted Display (HMD). Although the video content is displayed inside the HMD without any external exposure, we observe that user head movements are driven by the video content, which creates a unique side channel that does not exist in traditional 2D videos. By recording the user whose vision is blocked by the HMD via a malicious camera, an attacker can analyze the correlation between the user's head movements and the victim video to infer the video title.

To exploit this new vulnerability, we present INTRUDE, a system for identifying 360° videos from recordings of user head movements. INTRUDE is empowered by an HMD-based head movement estimation scheme to extract a head movement trace from the recording and a video saliency-based trace-fingerprint matching framework to infer the video title. Evaluation results show that INTRUDE achieves over 96% of accuracy for video identification and is robust under different recording environments. Moreover, INTRUDE maintains its effectiveness in the open-world identification scenario.

## 1 Introduction

Virtual Reality (VR) is a rapidly growing technology with a projected market value of $22 billion in 2025 [81]. Meta, one of the industry's leaders, has sold approximately 20 million Oculus VR headsets as of March 2023 [71, 84] and invested billions in its VR venture. Beyond the popularity in gaming and entertainment, VR head-mounted displays (HMDs) are transforming a wide range of industries, including military training, medical operation, and virtual conferencing [36, 65].

The immersive nature of VR HMDs has given rise to various prevalent applications, one of which is the 360° video. 360° videos are widely used in areas such as virtual tours, live concerts, and immersive storytelling [10, 60]. Unlike traditional 2D videos with limited viewing perspective on a 2D frame, 360° videos enable viewers wearing an HMD to freely explore the video content in all directions by simply moving their heads, as if they were physically present in the scene.

VR technology presents unique privacy challenges due to its collection and storage of extensive personal data [25, 39, 83, 94]. Researchers have identified several attacks aimed at inferring VR keystrokes [9, 47, 53, 58, 78, 93, 103]. In this paper, we unveil a new threat that identifies the titles of 360° videos being viewed inside the HMD. The leakage of video titles reveals significant information about users' personal interests, opinions, and even religions [95], which unintended third parties have exploited for pushing political agendas [2] or conducting blackmail [3]. Note that video identification is not necessarily related to user identification. For example, video identification attacks can target a group of users and reveal the political or religious interests of certain organizations or regions, without performing user identification. While side-channel attacks identifying 2D videos displayed on smartphones, PCs, and TVs have been previously studied [26, 31, 69, 75, 95], no research has demonstrated the feasibility of side-channel attacks on 360° videos displayed on a VR HMD.

**INTRUDE.** We present INTRUDE, a vIdeo ideNtification aTtack towaRd virtUal reality HMDs on 360° vidEos by leveraging a new *contactless* side channel – user head movements. While the VR content displayed inside the HMD remains inaccessible to external attackers, we discover that the head-movement-based interaction between the user and the HMD creates a side channel that is completely exposed to the public. As video content drives the user to move or fix his/her head [6], there is a subtle relationship between the user's head movement and the displayed 360° video, which can be exploited to infer video titles without direct access to the HMD. By taking advantage of the fact that the victim cannot see the physical world when using the HMD, the attacker can freely record the victim's head movements. After extracting the head movement from the camera recording and matching it with fingerprints of videos of interest, INTRUDE can infer the title of the playing 360° video.

1

Compared to prior camera-recording-based attacks on smartphones and PCs [20, 82, 101], INTRUDE is easier to launch and harder to detect as VR users are oblivious to their physical surroundings when wearing the HMD. The attackers have more flexibility in choosing when, where, and how to record the victim's head movement. Furthermore, since INTRUDE utilizes a new side-channel that does not exist in 2D video viewing, conventional countermeasures designed for 2D video identification, such as varying video encoding configuration for traffic analysis attacks and adjusting screen brightness for screen reflection attacks [75, 95], become ineffective.

**Challenges and Approaches.** Realizing INTRUDE requires addressing two new challenges in the context of 360° video identification on VR HMDs. *First*, it is non-trivial to extract the victim's head movement trace from the image pixels of the recording. Given that HMDs cover the majority of users' faces and users may look around in all directions during 360° video viewing, existing head pose estimators [16, 64, 72, 104] harnessing facial features for limited front-facing positions cannot be directly used. To overcome this challenge, we propose an HMD-based *head movement estimation* scheme to extract the full-spectrum head movement trace from the image pixels of the camera recording via a deep convolutional neural network (DCNN).

*Second*, the correlation between the video title and the head movement is implicit and subtle, making it challenging to fingerprint and identify the video. Fingerprinting a video directly by head movement traces is ineffective due to the fluctuation of these traces in response to changes in human conditions or environments. Consequently, treating the head movement trace as the fingerprint may lead to poor matching with the victim's actual trace, even if they embed similar head movement patterns. Hence, we propose to fingerprint a video by saliency maps as they can capture the stable features embedded within the fluctuating head movement traces across different trials with the same video. We design a *trace-fingerprint matching* framework to match the extracted head movement trace with the video saliency fingerprints via a multi-modality model, enabling video identification.

**Evaluations.** We validate INTRUDE through extensive evaluations that involve 31 users, 635 videos, and two HMDs. In the default indoor attack scenario, INTRUDE achieves the top-1,2,3 accuracy of 96%, 99%, 99% for 360° video identification. The robustness analysis under various recording environments shows that the attack maintains a stable video identification accuracy across different recording distances, lighting, backgrounds, and angles. Moreover, the risk of INTRUDE is validated in the *open-world* identification scenario.

**Ethical Consideration.** All experiments involving human subjects in this study have been approved by the IRB. We have only used INTRUDE to perform attacks on the datasets described in this paper. INTRUDE has never been used in other ways or released to other parties.

**Contributions.** The contributions of this paper include:

- We present the first 360° video identification attack on VR HMDs through a new *contactless* side channel – user head movements (§2 and §3). This threat has not been identified and is fundamentally different from previous attacks.

- We build new techniques to realize the attack, i.e., a head movement estimation scheme to obtain the victim's trace (§4) and a trace-fingerprint matching framework for video identification (§5). The cores of these two system components are two novel neural networks, one that can extract head orientations from recordings of 360° HMD viewing and the other one that can match different input modalities for video inference.

- We perform an extensive evaluation (§6 and §7) of the proposed attack system INTRUDE, including validation of system components and parameters, baseline comparison, robustness analysis, and open-world identification attacks.

## 2 Background and Motivation

### 2.1 360° Videos and Head Movement

360° video viewing is an important application in VR. To view a 360° video using an HMD, a user must first select the video of interest via the VR hand controller. Once the video playback starts, the VR engine of the video player decodes the video into a spherical scene wrapping around the user's head. A subset of the sphere to which the head is oriented, i.e., the viewport, is rendered and displayed on the HMD screen. This way, the user can simply move his/her head around to explore the 360° video in an immersive hands-free way.

It is well-known that the displayed 360° video drives the head movement of the user [27, 62]. This is because the human vision system tends to selectively pay attention to some salient objects or events in a visual scene instead of processing every single pixel. As a result, the user follows his/her visual attention and moves his/her head around when viewing the 360° video. The user may explore and glance through the content or fixate on a region of interest. Low-level content features of 360° videos, such as color, contrast, and textures, can affect visual attention and head movement. High-level content features such as human faces and notable objects are also critical to the head movement.

More importantly, such a correlation between the 360° video and the head movement are similar across users [17, 22, 50, 61]. In other words, different users would present similar head movement patterns when viewing the same video. Therefore, even though we have no knowledge about the victim, we can collect fingerprints for videos of interest and utilize this correlation to match the victim's head movement trace with all the fingerprints. By finding the most similar fingerprint, we can infer the video being viewed.
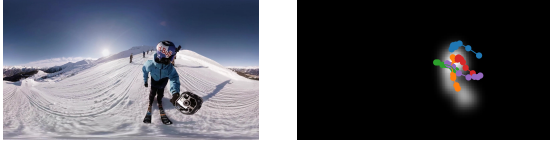
Figure 1: The head movement traces of different users (right figure, 4 colors) on the same video may be different in the time domain, but they do present a similar head movement pattern on regions of interest (right figure, brighter areas).
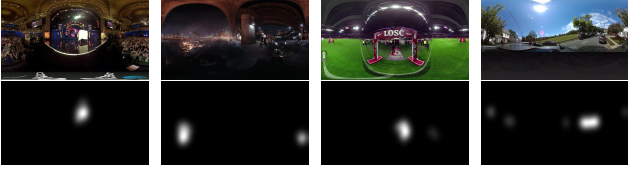


Figure 2: The saliency maps (bottom) uniquely identifies 360° videos (top) through regions of interest.

## 2.2 Video Fingerprint and Saliency Maps

Despite the intuitive principle, matching head movements with videos is challenging. Our study discovers that even though similar patterns exist in different head movement traces of the same video, it does not mean that these head movement time series are exactly the same. In fact, these head movement traces of the same video may differ from each other significantly in the time domain due to the time-varying physical and emotional states of the viewer and external factors of the viewing environment [6, 40]. For example, two viewers may look at the same object at different moments or view different parts of the same object at the same time. Figure 1 illustrates this fact via one-second head movement traces of four users (right) viewing a sample 360° video (left). We can see that different users' head movement traces (marked by four different colors) deviate from each other. Therefore, because of the randomness of head movement in the time domain, even if we can obtain the head movement trace for a video, it is infeasible to use it as the video fingerprint and directly match it with the victim's head movement.

Fortunately, we also observe that different users' head movements tend to focus on the same region, e.g., the bright area in Figure 1 (right). This is attributed to the fact that this region contains attentive content, i.e., the blueish upper body of the human, and thus drives all users' heads to move here. This phenomenon is consistent with previous findings in viewing behavior of 360° videos [17, 27, 62]. In essence, these salient regions of a 360° video are highly correlated to different users' head movement traces. They actually represent the similar head movement pattern across different users' traces. Hence, we propose to fingerprint 360° videos by their salient regions and then match the victim's head movement trace with the fingerprint (salient regions).



Figure 3: Sample attack scenarios – recording with a smartphone rear camera (left) and front camera (right).

The saliency of a video can be quantified by *saliency maps*. A saliency map is a 2D heat map that indicates the salient regions in a video frame. More salient regions where users stay for a long time are assigned higher saliency values and represented by brighter pixels in the saliency map, and vice versa for less salient regions that users glance through. The saliency map of an image can be generated by saliency detection, an established technology that has been widely used in surveillance [100], advertising [44], and content generation [37, 88]. As shown in Figure 2, the saliency maps in the bottom row signify salient regions in the frames of four 360° videos in the top row. These saliency maps contain distinctive information about a visual scene that uniquely identifies the 360° videos. For instance, the man on the stage in the left-most image is the most conspicuous. Note that saliency maps remain effective in fingerprinting video frames even when they contain multiple salient regions. For example, the second-from-the-left saliency map in Figure 2 has two salient regions. It is possible for different users to attend to either the left or right region, with the likelihood being proportional to the brightness of the salient region. However, regardless of which region the users focus on, their head movement traces would still be correlated with this double-region saliency map rather than a saliency map of another video frame that has no overlap with their head movements. Essentially, the saliency map essentially captures the collection of regions with frequent head movements.

As for a video, the saliency map changes over time for each frame, following the progress of the content. The saliency fingerprint of a video, i.e., a sequence of saliency maps, becomes even more distinct. Therefore, we can conclude that it is *feasible* to utilize saliency maps to fingerprint a 360° video.

## 3 Overview of INTRUDE

### 3.1 Threat Model

We consider an adversary with a library of target 360° videos whose viewers can be exploited for malicious purposes, e.g., pushing agendas or conducting blackmail. The attacker wants to know if a victim is viewing one of these videos and, if so, infer the specific video title. Existing studies [31,56,69,75,95] have confirmed the sensitivity and privacy associated with such video title information. In our attack scenario, the victim

is viewing a 360° video using an HMD in a public space (e.g., a library [4] or an airport [1]). Watching 360° videos publicly with VR headsets has become popular due to the low purchasing cost [99]. For example, 100% libraries from the Association of Research Libraries provided VR headsets for in-library usage [24]. There is also an increasing number of users watching 360° videos during a flight [68, 85]. At the same time, the attacker is recording videos of the victim's head movement using a filming device, such as a smartphone or a digital camera (examples shown in Figure 3). The captured videos are referred to as *recordings* and used to infer the 360° video viewed by the victim. This scenario is inspired by prior works where the attacker captures shoulder surfing recordings to infer the victim's unlock patterns on Android [43, 87, 97]. However, our case is more plausible because the victim's eyes are fully covered by the HMD and thus the victim is less vigilant about the proximity. The attacker may choose to place a hidden camera to record the victim's head movement to further reduce suspicion.

We assume the attacker has physical access to the location where the victim is viewing 360° videos and can capture the recording. Also, we assume that the attacker knows the type of HMD the victim is using, which is trivial to obtain since the attacker can visually recognize it. The HMD type information will be used to tune the system components of INTRUDE. We do *not* assume any other capabilities for the attacker. For example, the attacker cannot lure the victim into downloading and installing a malicious App [78], and the attacker cannot sniff the network traffic coming from the victim's device.

## 3.2 Challenges and Key Components

To identify 360° videos from recordings of victim head movement, we need to tackle the following challenges.

- **Extracting Head Movement Traces in VR**. Head movement tracking in VR presents a distinct challenge compared to traditional head pose estimation in 2D scenarios [59]. VR head movement encompasses omnidirectional viewing directions and the user's face is covered by the HMD. As a result, prior head pose estimators [16, 42, 72] assuming that head poses and facial features are explicitly displayed in front of the camera in a limited range of orientations would be ineffective.

- **Identifying Videos by Trace-Fingerprint Matching**. Another challenge is that the video title is hidden inside the spatial-temporal correlation between two different data sources – the head movement trace and the video fingerprint. Unlike previous video identification attacks that utilize a univariate time series, e.g., network traffic [31] or light effusions [95], to fingerprint a video and match it to the victim's corresponding time series, fingerprinting a video by fluctuating head movement traces is ineffective for trace-fingerprint matching and video identification.
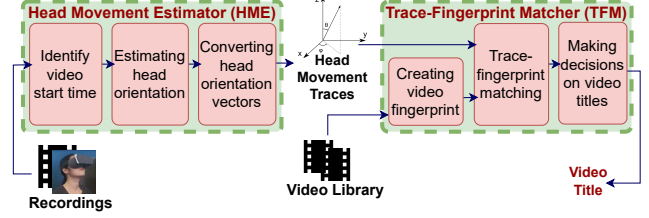


Figure 4: Overview of the proposed INTRUDE system.

INTRUDE overcomes these challenges through its two components as shown in Figure 4. First, the head movement recording of the victim while viewing an unknown video is passed to the HMD-based head movement estimation scheme (HME). This scheme accurately identifies and extracts visual head movements during the omnidirectional viewing. The extracted movements are then converted into a head movement trace, a time series of coordinates tracking the head orientations relative to the VR coordinate system throughout the viewing session. Second, the head movement traces are fed to the video saliency based trace-fingerprint matching framework (TFM). The trace-fingerprint matching model matches the victim's trace with each video fingerprint in the library of target videos and produces a video identification decision.

## 4 HMD-based Head Movement Estimation (HME)

This section presents INTRUDE's head movement estimation scheme. INTRUDE identifies the start of video playback from the recording, extracts head orientations using the proposed DCNN, and generates the head movement trace after representing the head orientations in the VR coordinate system.

### 4.1 Identifying the Start of Video Playback

A prerequisite to launching INTRUDE is determining if the victim is viewing a 360° video or engaging in other VR activities. Filtering the recording to obtain the 360° video viewing session is necessary because the recording may contain redundant information generated when the victim is using other VR Apps. If the victim is viewing a 360° video, it is also critical to identify the moment when the 360° video starts playing on the HMD screen. After this moment, the victims will only move their heads based on the video content. Hence, this moment should be identified as the start of the head movement trace for launching the attack.

This challenge can be overcome by directly inspecting the recording due to the unique interaction pattern in 360° video viewing sessions. Initiating the playback requires the user to navigate to the desired video and "click" it for playback. Despite the differences in user interfaces in different HMDs, the navigation and selection of the video are always executed by a set of head and hand actions, ending with a press on

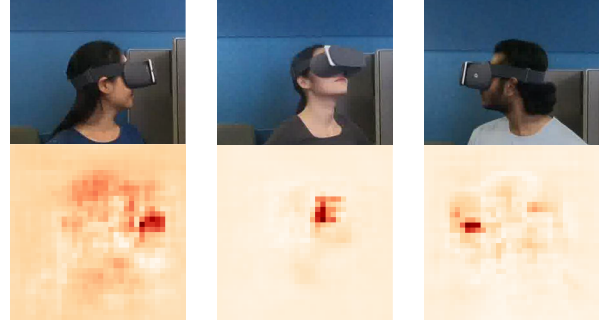Figure 5: The head orientation estimation model.



Figure 6: The head orientation estimation model focuses on learning HMD features (dark red areas in the bottom row) and can generalize well to victims not presented in the training.

the hand controller for starting the playback. After this final clicking, the viewing session starts and the victim only moves his/her head to explore the 360° video without further operation on the controller. This head-only, hand-free interaction design ensures immersive user experience [57, 86], differing from other VR Apps that always require both hand controller and head movement. By recognizing this unique sequence of user actions through visuals and sounds in the camera recording, we can infer that the victim is viewing a 360° video and pinpoint the moment of the final controller clicking as the playback start time. INTRUDE extracts all frames from this starting moment for a recording length of $T$ seconds.

We then locate and crop the regions that are relevant to the attack in the filtered recording. This can be easily performed by off-the-shelf video editing software because victims interact with 360° videos only through head movement and their body positions remain stable during the viewing. Consequently, the recording becomes a smaller version only containing the victim's head, shoulders, and HMD. This preprocessing reduces the input size and thus the computation cost. It also eliminates noise in the background pixels, boosting the head movement estimation performance. After cropping, the recording is split into individual video frames before being sent to the head orientation estimation model.

## 4.2 Estimating Head Orientation

This section details our head orientation estimation model's design, which analyzes the processed frames in the recording. The core of the head movement estimation (HME) component, the head orientation estimation model, transforms visual data from the *recording* into numerical head orientation data. Previous estimation methods [16, 64, 72, 104] extracted head orientations from general non-HMD head pose images. By assuming that users mostly look straight ahead into the camera, they relied on facial landmarks to derive head orientations. However, these models fail when users' faces are obscured by HMDs and users turn their heads away from the camera's filming direction during 360° video viewing. In INTRUDE, our head orientation estimation model is specifically designed to estimate HMD users' head orientations without the limited range of head movements assumed by prior works. It focuses on learning the visual features of HMDs, which are more discernible than facial landmarks, particularly when users deviate 45° or more from the camera's filming direction. Therefore, our model can support the estimation of the full spectrum of head movements in 360° video viewing.

**Estimation Model Design.** The head orientation estimation model architecture, depicted in Figure 5, is a deep convolutional neural network (DCNN) comprising a base network and a decision layer. The base network explores pixel values of the input frame. It extracts and learns the representation of meaningful features for the head orientation estimation. The decision layer estimates the head orientation based on the output of the base network. Inspired by the ResNet-50 model [34], our base network consists of 48 stacked convolutional layers and an adaptive average pooling layer. The ability to learn complicated visual patterns of HMDs through this highly non-linear structure is our key advantage in head orientation estimation compared with previous approaches using hand-crafted features. The adaptive average pooling layer receives a set of feature maps from the base network. It then calculates the mean values of each feature map and returns a flattened vector. By fixing the length of the flattened vector, our design can handle different sizes of input frames. Finally, the dense decision layer maps the flattened vector to a *head orientation vector*, a 3D vector representing the head orientation in 3D spaces.

**Model Training and Operation.** We use Mean Squared Error (MSE) loss to guide the learning of the head orientation estimation model. It calculates the mean squared distances between the estimated and the ground truth head orientation vectors and forces the model to push its estimation to the ground truth. Formally, MSE is defined as:

$$\mathcal{L} = \frac{1}{3}((x - \hat{x})^2 + (y - \hat{y})^2 + (z - \hat{z})^2) \qquad (1)$$

where $v = (x, y, z)$ is the estimated head orientation vector extracted from input frame $f$, and $\hat{v} = (\hat{x}, \hat{y}, \hat{z})$ is the ground truth head orientation associated with $f$.

Since the above head orientation estimation model is not a personalized model designed only for a particular victim, attackers can obtain training data offline by rotating their heads to different angles with respect to the camera's filming direction, just as normal users. They can capture the camera recording of their own head movements by using the same

type of HMD as the victim. Meanwhile, the ground truth head orientations can be collected by the HMD sensors. Once the model is trained, it can then be used to estimate the head orientation of the victim who is not in the training set, i.e., a testing user. Figure 6 illustrates that the trained model focuses on learning the features of HMDs and can be well generalized to new testing users. Through the occlusion method [102], the heat maps at the bottom row highlight the regions affecting the estimation performance for the testing frames in the top row. These regions contain HMDs (darker red) rather than heads, ears, or jaws (lighter red) and are consistent across victims wearing similar HMDs. After applying the model to each video frame in the recording, a sequence of head orientation vectors is eventually generated.

## 4.3 Converting Head Orientation Vectors

The proposed estimation model generates a sequence of head orientation vectors, but these cannot directly form a head movement trace because they are referenced to the recording camera's filming direction, i.e., the *camera-based coordinate system*. This differs from the *VR coordinate system* used in VR Operating Systems (OS) for tracking head movements and rendering VR content. As a result, an extracted vector may point to a totally different direction in the VR OS. Hence, INTRUDE must convert these vectors to the VR coordinate system before inferring a victim's regions of interest and video title.

Figure 7 illustrates the difference between the two coordinate systems. It can be seen that both systems have an axis pointing straight up (the *z* and yaw axes), but their reference axes are distinct. The reference axis in the camera-based coordinate system (the *x* axis) aligns with the camera's filming direction and can be visually identified in the camera recording. On the other hand, the reference axis of the VR coordinate system (the roll axis) is represented in the VR virtual world and cannot be observed in the physical world by the attackers.

To locate the VR coordinate system, we first identify the roll axis by leveraging the fact that the VR coordinate system is not fixed in the physical world but is always reset by the VR OS at the beginning of a $360°$ video viewing session. Upon reset, the roll axis is configured to overlap with the projection of the head orientation vector onto the *x*-*y* plane. By extracting the head orientation vector at the reset, i.e., at the starting moment of the video playback, we locate the roll axis in reference to the camera-based coordinate system.

INTRUDE can then convert a head orientation vector for it to be represented in the VR coordinate system. As shown in Figure 7, the two systems already have two axes aligned, namely the yaw and *z* axes. Thus, the conversion becomes the angular rotation along the *z* axis for the *x* and roll axes to align. As discussed above, the offset angle between the *x* and roll axes can be derived when locating the roll axis. Without loss of generality, we denote the first frame at the beginning
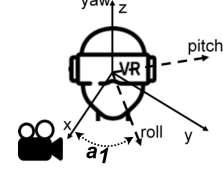


Figure 7: Difference between the camera-based coordinate system (solid lines) and the VR coordinate system (dash lines).

of the viewing session as $f_1$ and the head orientation vector in the camera-based coordinate system extracted from $f_1$ as $v_1$, $(x_1, y_1, z_1)$ being the coordinate of $v_1$. Let $P$ be the plane formed by the intersection of the *x* and *y* axes. Then the offset angle $a_1$ between the *x* axis and roll axis (the projection of $v_1$ onto $P$) can be expressed as $a_1 = arctan(\frac{y_1}{x_1})$. Given the offset angle $a_1$, the mapping $q_1$ is the quaternion rotation operation [15] rotating a 3D vector for $-a_1°$ along the yaw axis. Applying $q_1$ to a head orientation vector $v$ in the camera-based coordinates produces the corresponding representation in the VR coordinate system $u$.

## 5 Video Saliency-based Trace-Fingerprint Matching (TFM)

This section details INTRUDE's approach to inferring video titles through trace-fingerprint matching. INTRUDE generates video saliency-based fingerprints for $360°$ videos targeted by the attacker, and matches the victim's head movement trace obtained from the HME component with each fingerprint to infer the video being viewed.

### 5.1 Creating Fingerprints for $360°$ Videos

An essential offline task for attackers before making the video identification decision is to build a library of videos of interest and create fingerprints for these videos. These videos of interest can be carefully selected based on the malicious purpose. For example, in the case of blackmailing, sensitive videos that the victim is not supposed to view can be put in the library. As discussed in §2.2, a head movement trace of a single user is insufficient to fingerprint a $360°$ video because the head movement time series generally fluctuates. However, users tend to agree on the regions of interest in $360°$ videos. This fact drives the similar head movement patterns across users that center around salient regions of a video. Therefore, we propose to utilize video saliency maps to fingerprint a $360°$ video.

A saliency map can be derived offline by collecting and processing head movement traces of a group of users, but it would be too expensive to collect subjective human data for every $360°$ video targeted by the attackers. INTRUDE instead uses an established saliency detection model to generate the saliency map. Thanks to the recent development of deep learning, salient regions of human users on 2D videos and $360°$

videos can be reliably and accurately detected [38, 62]. These state-of-the-art saliency detectors have been trained on a large amount of data and can identify the salient regions of new unseen videos. To fingerprint a video, INTRUDE performs saliency detection on all its frames through a state-of-the-art model. The output of the model, a sequence of saliency maps, is a spatial-temporal video signature reflecting the progress of human attention through time and is stored as the video fingerprint. This procedure is repeated for all library videos.

## 5.2 Matching Victim Trace with Fingerprints

Unlike the majority of previous works that only processed time series to identify video titles [47, 56, 75, 95], INTRUDE needs to handle input of different modalities to infer the video title. The victim's head movement trace is a time series of head orientation vectors, whereas the fingerprint of a $360°$ video is a sequence of 2D saliency maps. INTRUDE's trace-fingerprint matching model must establish the spatial correlation between the victim's head movement and the video's salient regions, as well as comprehend the head movement patterns in relation to the progression of the video's salient regions. Such tasks are challenging for traditional machine learning models hand-crafted by domain knowledge [31, 69, 75]. To overcome this, we propose a DCNN to capture the interactive patterns between the two input modalities, effectively matching the victim's head movement trace to a library of video fingerprints for video identification.

**Input Preprocessing.** Since the victim's head movement trace and the video fingerprint are two different modalities, they cannot be analyzed simultaneously. We propose to first transform the head movement trace into a sequence of 2D head orientation maps. Let $\mathcal{V}$ be the head movement trace and $\mathcal{M}$ be the video fingerprint, where $\mathcal{V} = v_1, .., v_N$ is a list of 3D head orientation vectors in the VR coordinate system and $\mathcal{M} = m_1, ..., m_N$ is a list of saliency maps. Each saliency map $m_i$ is a 2D heat map of size $W \times H$. We create the head orientation maps through equirectangular projection [67], a prevalent approach for projecting 3D spherical data onto a 2D plane. We adopt this method because modern $360°$ video saliency maps are marked on equirectangular 360 frames [61, 62]. This projection aligns both input modalities. For a 3D head orientation vector $v_i$, we find the azimuth and altitude angles, $\theta_i$ and $\phi_i$. Then, the projected point $v'_i = (w_i, h_i)$ on the equirectangular of size $W \times H$ can be derived using the following formulas [61],

$$w_i = (\frac{\theta_i}{360})W, \ \ h_i = (\frac{1 - sin(\phi_i)}{2})H. \quad (2)$$

After this conversion, we obtain a sequence of head orientation maps of size $W \times H$ where the head orientation is indicated by $\mathcal{V}' = (v'_1, ..., v'_N)$.

We then combine the transformed head orientation maps with the video fingerprint into a unified input. Given that both head orientation maps and saliency maps are associated with
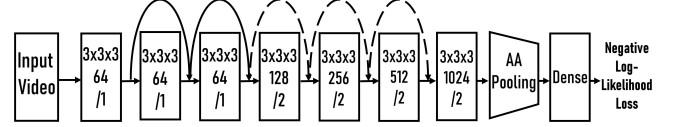


Figure 8: The DCNN architecture of the video saliency based trace-fingerprint matching model.

timestamps that signal their temporal position in the video, we combine the corresponding maps of the same timestamp. The result is a sequence of pairs of head orientation maps and saliency maps that can be processed by our DCNN. The corresponding regions in both head orientation maps and saliency maps are aligned. Consequently, the DCNN better learns the spatial relevance between the victim's head orientation and the video's salient regions.

Due to the computational intensity of the DCNN, processing the entire sequence of head orientation and saliency maps can be inefficient. Yet, feeding a sparse number of maps to the DCNN may fail to present a continuous context of how the victim's head movement interacts with the video saliency. To balance efficiency and performance, we introduce $\tau$, a system parameter defining the time interval between two samples in the sequence. By adjusting $\tau$, INTRUDE can flexibly control the sampling interval of input data for the trace-fingerprint matching model. For instance, setting $\tau = 1$ directs our system to sample the sequence of head orientation maps and saliency maps every second.

**Matching Model Design.** The matching between the preprocessed head movement trace (head orientation maps) and video fingerprints (saliency maps) can be regarded as a mapping problem, where the two input channels are mapped to a confidence score that indicates the correlation between them. The core of the trace-fingerprint matching model is a DCNN, as depicted in Figure 8, composed of seven stacked 3D convolutional layers. Unlike traditional 2D convolutional layers, the 3D kernel convolves over several input maps at a time to learn both the localized spatial information in the head orientation maps and saliency maps, and the short-term temporal relationship between these maps. Stacking several 3D layers expands the neuron receptive fields, enabling the model to learn long-term information throughout the video. This design is consistent with the spatial-temporal structure of the two input modalities to match.

For each convolutional layer, there are three associated parameters, the kernel size, the number of kernels, and the stride. The kernel size of $3 \times 3 \times 3$ is selected for all convolutional layers to reduce the computational cost. Kernels having the stride value of 2 (denoted by $/2$) reduce the size of the intermediate maps by half. We use residual connections [34], shown as U-shape arrows in Figure 8, to connect the output of the previous layers with the output of the next layers. The connections facilitate the backward flow of gradients from

7

the output back to the input, mitigating the impact of the gradient vanishing. Similar to the head orientation estimation model, the end of the DCNN is the adaptive pooling layer [33] converting the filtered maps of arbitrary sizes into a fixed flattened vector. The decision network is a dense layer of 1024 neurons mapping the flattened vector to a single confidence score in the range of $(0, 1)$.

To measure the difference between the output confidence score of the model and the ground truth, we use the Negative Log-Likelihood loss,

$$L(n, \hat{n}) = -\hat{n} log(n) \qquad (3)$$

where $n$ is the output of the proposed model and $\hat{n}$ is the ground truth. If the head movement trace and the video saliency maps for matching indeed come from the same video, $\hat{n}$ has the value of 1. Otherwise, $\hat{n}$ has the value of 0. Since $\hat{n} \in \{0, 1\}$, the value of $L(n, \hat{n})$ can be either 0 or $-log(n)$.

## 5.3 Making Video Identification Decisions

As the trace-fingerprint matching model is not victim- or HMD-dependent, attackers can collect offline data of head movement traces and video fingerprints to train the proposed DCNN. The head movement traces can be collected by HMD sensors when attackers view target 360° videos and/or by directly using public head movement traces of popular 360° videos targeted by the attackers. The video fingerprints, i.e., the saliency maps, can be created by applying a high-performance saliency detector on the respective 360° videos. Once the model is trained, it can be used for video identification. To perform the attack, the victim's head movement trace is paired with each video fingerprint in the library and processed through the trace-fingerprint matching model. This outputs confidence scores that indicate the match between each video and the victim's head movement. INTRUDE then produces the top-$k$ candidates by selecting the $k$ videos with the highest confidence scores.

## 6 Experimental Methodology

### 6.1 Settings

**Apparatus.** The INTRUDE prototype comprises an HMD device, a GoPro HERO6 camera, and an offline processing desktop running on Ubuntu 18.04 equipped with a GeForce GTX 1080 Ti GPU. The GoPro6, which records videos in 1080p with a 12MP lens, is inferior in recording capability compared to modern smartphones such as iPhone 14 that is equipped with a 48MP lens and can produce 4K videos. We tested two cordless HMDs commonly used in VR studies [19, 30, 80], Google DayDream and Google Cardboard. Both have 100° field of view similar to most modern VR headsets. The camera was situated 1.5 meters in front of the user at head level, and recorded at 1080p resolution. We took recordings in

a well-lit office ($\sim$1000 lux as per [5]) at a sampling interval of $\tau = 0.8$. We extracted 60 seconds from each recording for our experiments. Half of them were extracted from the first 60 seconds, while the other half were segmented from a random point of the recordings. This default setting is maintained unless otherwise stated when evaluating specific factors.

**Subjects.** We recruited 31 subjects (15 males and 16 females, aged from 19 to 36) for the user studies. Subjects are university students and people recruited through public channels, among which 18 wear glasses, 5 have no experience in VR, and 10 have little experience. All subjects gave consent to our IRB-approved study, which allows us to record human responses for VR system evaluation. Subjects stood in an open space while viewing the provided content. All 31 subjects participated in model training (§6.2) while 17 subjects were randomly chosen to play the 'victim' role during the data collection phase (§6.3). The decision to work with 17 victims is attributed to the significant time commitment of playing a victim, aligning with the participant count in similar studies [78, 93]. For training the HME, recruits were asked to follow specific instructions to collect head orientations. When playing the role of victims, users were instructed to engage freely with the 360° videos without any additional guidance. All participants were kept unaware of the study's purpose.

**Video Library.** We collected a library composed of the top 635 popular 360° videos with the most views on YouTube, covering a diverse array of genres such as documentaries, music, sports, animations, wildlife, vlogs, travel, and gaming. Each genre contains a minimum of 30 videos and each video receives at least 40,000 views. We obtain the fingerprints of these videos by applying a state-of-the-art 360° video saliency detector [61].

### 6.2 Offline Model Training

**The HME.** To train the head orientation estimation model in the HME (§4.2), we first collected a dataset of camera recordings that captures different head orientations of the 31 subjects. Subjects were asked to view a target virtual object displayed on the HMD screen at different angles with respect to the camera filming direction. The head orientations during this one-second viewing session ranged from -179° to 179° in yaw and from -55° to 55° in pitch. Note that subjects were not instructed to freely explore 360 videos because that may not cover the diverse head orientations needed for HME training.

Each subject repeated the viewing session 15 times for one HMD device. We then annotated the head orientation for each frame of the recordings as the ground truth. A total of $\sim$200,000 annotated head orientation frames were collected to train the model. To mitigate model overfitting and increase the variety of patterns the model can learn, we applied data augmentation to the training data. We randomly picked a method from horizontal flipping, adding random gray boxes,

modifying brightness and contrast, and zooming in and out at the center, and then iterated through all frames [76]. The total time to train the head orientation estimation model was approximately 5 hours. Details on creating the dataset can be found in the Appendix.

**The TFM.** As discussed in §5.3, the training process of the matching model uses head movement traces when viewing videos and fingerprints of these videos. An inherent advantage of this design is that TFM training does not rely on camera recordings of head movements, making it independent of the training and testing performance of the HME. This enables attackers to leverage publicly available head movement traces for TFM training, which offers more accurate data than the HME output and reduces the data collection efforts. We implemented this design by utilizing three existing datasets of head movement traces for 360° videos [22, 50, 91]. These traces were collected by built-in HMD sensors for 24 videos, each with head movement traces from at least 48 users in a free-viewing context. In addition, we obtained the fingerprints of these videos by applying the aforementioned saliency detector [61]. By combining a head movement trace and a saliency fingerprint as one sample, we created a training set of 1,152 positive samples and 23,392 negative samples, where *positive* means the trace matches with the video title and *negative* indicates no match.

To mitigate the overfitting, we applied shifting augmentation [76] on head orientation maps, where the pixels indicating head orientations were randomly shifted up, down, left, or right to simulate the head orientation estimation errors. We leveraged an upsampling technique to mitigate the impacts of imbalanced data. The total training time was about 16 hours.

## 6.3 Online Attacks

Once the models in the HME and the TFM are trained, we launch INTRUDE to execute attacks and evaluate the performance of video identification. We conducted *cross validation*, wherein each round one user played the role of a victim and his/her data was excluded from the training set to ensure no overlap. We simulated the scenario where the victim viewed some videos in the target library, while the attacker aimed to determine which videos were viewed. We randomly selected 22 videos from the 635-video library for the victim to view, resulting in 22 camera recordings of head movements. For each recording, INTRUDE inferred the top-*k* video titles out of the 635 potential options. We repeated this process for 17 victims and report the average results.

## 7 Evaluation Results

In this section, we present the extensive evaluation results of INTRUDE when the online attacks are launched. We start with validating the head movement estimation (HME) (§7.1). Next, we focus on evaluating the video identification attacks,
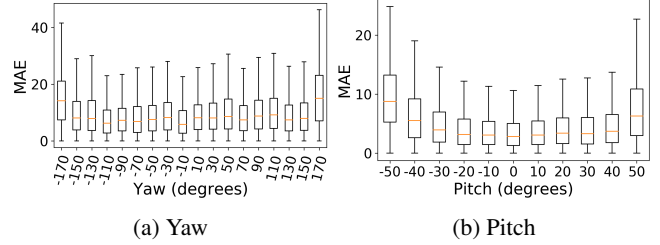


(a) Yaw          (b) Pitch

Figure 9: Estimation errors of the HME in yaw and pitch.

including analysis of system parameters (§7.2), comparison with baseline approaches (§7.3), and robustness analysis under various environmental factors (§7.4). Finally, we evaluate the open-world identification (§7.5).

### 7.1 Validation of the HME

**HME Effectiveness.** We first validate the effectiveness of the HME because it extracts head movement traces from camera recordings and affects the performance of INTRUDE. We use Mean Absolute Error (MAE) to measure the absolute angles between the estimated and ground truth head orientation vectors. We focus on the yaw and pitch angles because the roll angle is not needed to locate where the victim is looking at the 360° video. It only indicates how the head tilts (not turns).

Figure 9 shows the boxplot of MAE (median values marked in red) for different head orientations categorized by their respective yaw and pitch coordinates. The average MAE between estimated and ground truth head orientations are 8.8° and 4.3° for the yaw and pitch coordinates, respectively. These errors are comparable to those obtained from generic head estimators for non-HMD limited-range head orientations [16, 42, 72]. Note that the MAE for the yaw coordinate is higher than the pitch coordinate. This is due to the yaw drift effect [28] where gyroscope sensors incur noises in the yaw coordinate of the ground truth and lead to biases in the yaw estimation. We also observe that the HME performs well in a wide range of yaw and pitch angles supporting 360° video viewing. The performance degrades at extreme angles where fewer samples were collected for model learning. However, since users rarely move to those places, the overall effectiveness of the HME is not affected.

**Head Tracking on Different Devices.** Next, we test HME performance with two Google Cardboard and Daydream, two popular VR headsets with different shapes, sizes, materials, and colors. We aim to explore how different HMDs in the camera recording would affect the attack. The average MAE over all head orientations, the MAE of front head orientations looking toward the camera, and the MAE of back orientations turning back away from the camera, are shown in Figure 10. In general, the MAE of head orientation estimation for Daydream and Google Cardboard are similar, at an average of
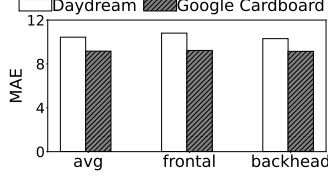
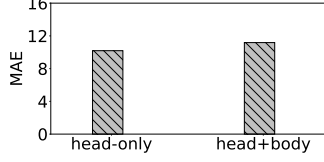Figure 10: Estimation errors of the HME for different VR headsets.



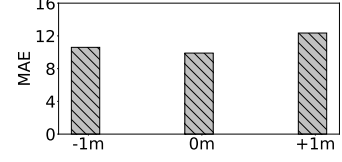Figure 11: Impact of body movements on HME tracking.



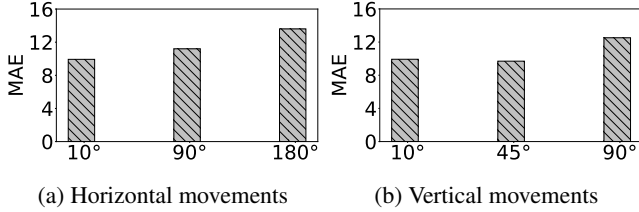Figure 12: Estimation errors of the HME at different camera heights.



(a) Horizontal movements      (b) Vertical movements

Figure 13: Impact of deliberate head movements on HME tracking.

$10.4°$ and $9.16°$ respectively, suggesting INTRUDE's capability of tracking head movements of victims wearing different HMDs in the recordings.

**Head Tracking with Body Movements.** To explore the impacts of body movements, we retrained the HME using two recording cropping methods, one only containing head movements (head-only) and the other containing both body and head movements (head+body). We present the results of MAE in Figure 11. As shown in the figure, the MAE errors of the head+body setup were 10% higher than those of the head-only setup since extracting HME-relevant features becomes more challenging when both head and body movements are present. This observation implies background body movements should be cropped out, as used in our default setup.

**Head Tracking with Deliberate Movements.** Deliberate head movements may disrupt head tracking if they cause the HMD to turn away from the camera, thereby obscuring visual features vital for tracking. We explored the effects of two types of deliberate head movements, vertical angular shifts at 10, 90, and 180 degrees, and horizontal angular shifts at 10, 45, and 90 degrees. For each angle, we collected 20 viewing traces from four users. For each angular shift, users performed five deliberate movements at random moments during a 360 video viewing session and move back to their original positions. The total duration of these movements accounts for about 10% of a 60-second viewing session without affecting the regular viewing experience. The average MAE scores are reported in Figure 13. We observe minimal performance degradation for moderate movements, e.g., a vertical shift at 10 or 45 degree and a horizontal shift at 10 or 90 degree, while large deliberate movements at a 180-degree horizontal shift or 90-degree vertical shift decrease tracking performance noticeably. However, these extreme movements rarely occur in practice [22] as they require considerable muscle effort and disconnect users from the focused virtual content for a few seconds.

**Head Tracking At Different Camera Heights.** We now investigate the performance of HME when the camera is placed at different heights, i.e., the camera being one meter higher than (+1m), at the same height as (0m), and one meter lower than the user's head (-1m). For each height, we recorded 20 traces collected from four users. The camera was positioned 3 meters away from the user and the average MAE scores are reported in Figure 12. The MAE errors increase at -1m and +1m because positioning cameras at a height exaggerate the scale of the head movement with respect to the camera. Slightly looking down at 0m appears to look straight down when the camera is placed at +1m. These extreme poses that are challenging to estimate occur more frequently when the camera is placed at -1m and +1m, leading to more errors.

## 7.2 Analysis of System Parameters

We now proceed to evaluate the identification performance. If one of the top-$k$ inference results is the actual video being viewed, the attack is considered successful. We report the top-$k$ identification accuracy (i.e., attack success rate) across all victim recordings. In this section, we analyze two system parameters, the recording length and the sampling interval.

**Recording Length.** Longer recordings provide more user head movement data and more clues about the video content being viewed. However, recording the victim for too long increases the risk of exposing the attack. Thus, it is necessary to understand the exact impact of the recording length $T$, the system parameter defined in §4.1. We measure the identification performance when the recording length varies from 10 s to 60 s. We cap the recording length at 1 minute since today's 360° videos typically have several minutes of duration [8]. The results are shown in Figure 14a. As expected, the identification accuracy increases as the recording length increases. The top-1,2,3 accuracy of INTRUDE steadily increases to 96%, 99%, and 99%, respectively, when $T = 60$. This result suggests a small trade-off between the stealthiness and performance of the attack. In case a long recording is not possible, INTRUDE can utilize a shorter recording, e.g., 20 s, to achieve a lower top-1,2,3 accuracy of 85%, 91%, and 93%, respectively.

**Sampling Interval.** To minimize the input data while ensuring the attack performance, we study how the sampling interval of the input would affect the top-1,2,3 identification accuracy. We set the $\tau$ parameter discussed in §5.2 as 0.5, 0.8,
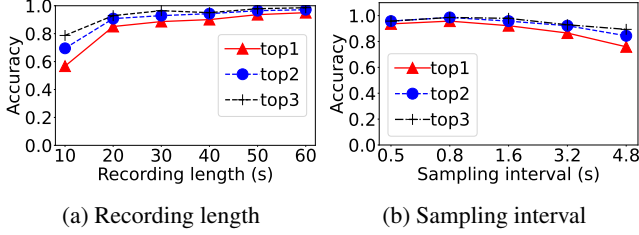
(a) Recording length      (b) Sampling interval

Figure 14: Video identification accuracy versus system parameters.



Figure 15: Benchmarking video identification accuracy by comparing INTRUDE with three baselines.

1.6, 3.2, and 4.8 seconds and repeat the attack to all recordings of victims. As shown in Figure 14b, INTRUDE achieves outstanding performance when $\tau$ is 0.5 or 0.8 seconds. For example, the top-1,2,3 accuracy reaches 96%, 99%, 99%, respectively, when $\tau = 0.8$. It can also be seen that the accuracy decreases as $\tau$ increases. The reason is that a greater $\tau$ makes the input information sparse, which imposes a penalty on the identification accuracy. On the other hand, minimizing the interval (e.g., 0.5 s) does not necessarily improve the accuracy because it could incur duplicated data samples that do not contribute meaningful information. Thus, we will use $\tau = 0.8$ for the rest evaluations.

## 7.3 Baseline Comparison

We also benchmark the performance of INTRUDE against the following baselines.

- **truth**: Instead of using head movement traces extracted by the HME, the trace-fingerprint matching model in *truth* directly uses the ground truth head movement traces collected by HMD sensors. Since *truth* is not impacted by estimation errors of the HME, its performance can be served as the upper bound.

- **htrace**: Instead of creating saliency fingerprints, a single head movement trace of a non-victim user is used as the video fingerprint in *htrace*. As head movement traces fluctuate significantly and are not ideal for video fingerprinting (see §2.2), *htrace*'s performance serves as the lower bound.

- **shallow**: To obtain insight into the importance of the design of the trace-fingerprint matching model, we create *shallow*. It is a variant of INTRUDE using a reduced version of the trace-fingerprint matching model, where the first two 3D convolutional layers are removed.

We repeat the experiments in the default setup for each baseline and report the results in Figure 15. As expected, *truth* has the highest accuracy since it utilizes the ground truth as the input for trace-fingerprint matching and video inference. The performance of INTRUDE closely approaches the upper bound, achieving a slightly lower yet comparable accuracy to *truth*. Fewer convolution layers in *shallow* reduces the top-1 accuracy by 1.4% because it affects the analysis of the fingerprints and head movement traces and therefore degrades the video identification performance. As for *htrace*, the
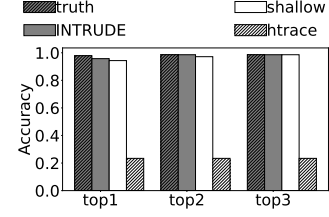
identification accuracy is the lowest, indicating that saliency fingerprints are essential features to assist the video matching.

## 7.4 Robustness Analysis

We now evaluate the robustness of INTRUDE under different environmental factors by repeating the attacks on all the victims. In each robustness study, one environment factor was different from the default setup, e.g., light condition. We captured these additional recordings of victims and launched the attack toward them. Since *truth* uses ground-truth head movement traces rather than extracted traces from recordings, it is not affected by the recording environment and thus is excluded in this evaluation. We report the top-1,2,3 accuracy of INTRUDE and top-1 accuracy of other baselines.

**Distances.** To measure the impact of the distance between the victim and the recording camera, we place the camera at 1.5 m, 3.0 m, and 4.5 m, away from the victim, respectively. The results in Figure 16a show that the performance of INTRUDE is stable as the distance varies. The attack distance only affects the clarity and visual details of the recordings and the extracted head movement trace. As long as the number of captured pixels for the head is sufficient, a long distance does not affect the identification accuracy. The distances in our evaluations are greater than or equal to those in prior recording-based attacks [47, 97]. To further extend the distance, we can utilize an advanced camera to capture 2K or 4K videos. The accuracy of *shallow* decreases a little as the distance increases because it is a weaker model with a low capacity to learn spatiotemporal patterns in the head movements and video fingerprints.

**Light Conditions.** The proposed attack may be launched in various light conditions depending on when the attacker finds the victim vulnerable. If the camera is placed in a dim environment, the recording could be noisy and might contain indiscernible regions [18]. Extracted head movement traces can then become erroneous, eventually affecting the attack performance. We launch INTRUDE on recordings captured in the bright indoor condition (~1000 lux) and the dark indoor condition (~50 lux) according to the lighting definitions in [5]. The bright condition is equivalent to a lit office while the dark condition is similar to a dark corridor or an unlit parking lot.
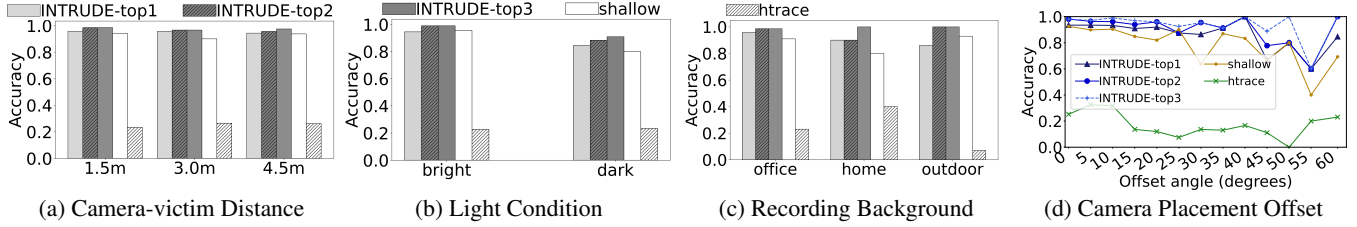
Figure 16: Robustness analysis of INTRUDE under different environmental factors.

As shown in Figure 16b, there is a small performance drop due to the effect of the dim environment. The top-1 accuracy of INTRUDE decreases from 96% to 84% while the top-1 accuracy of *shallow* drops from 92% down to 79%. However, we believe INTRUDE still poses a threat since the top-2,3 accuracy in the dim environment can reach 90% and 93%.

**Background Variations.** The location of the attack affects the background scenes in the recording. These visual backgrounds of the victim are analyzed to estimate head orientations and thus have an impact on the performance of INTRUDE. Therefore, we further investigate INTRUDE in different background scenes. In addition to the open office used in the default setup that includes gray cubicles and blue walls, we conduct the attack in two additional backgrounds. One is the home living room where the background scene consists of wood furniture, brown doors, and dark orange walls. The other is the outdoor background that is a typical garden with lawns and plants. Within each background category, we make the exact scene different for each victim to ensure the diversity of the recordings under evaluation. For example, different plants may appear behind a subject in the outdoor case. The results are shown in Figure 16c. Despite the slight variation, INTRUDE achieves a satisfactory performance across different backgrounds. The average top-1,2,3 accuracy across three backgrounds is 91%, 96%, and 100%, respectively. The support of different backgrounds stems from the HME design. As illustrated in §4.2, the HME learns to focus on the HMD device and ignore the background noises.

**Camera Placement Offset.** In the default setup, the camera is placed in front of the user along the user's front-facing direction, i.e., there is no offset angle between the filming direction and the front-facing direction (see §4.3). In this section, we place the camera in different locations to vary the offset angle between the filming direction and the head's front-facing direction. When the offset angle increases, the camera is gradually moved to the side of the victim. The results in Figure 16d show that there is a slight decrease of top-1 accuracy when the offset angle increases, e.g., around 85% for INTRUDE at the offset angle of 60°. This is due to the increased errors in head orientation estimation when the face and the HMD rotate away from the filming direction (see §7.1). Despite the decrease in top-1 accuracy, we note that INTRUDE still performs well, with a top-2 and top-3 accuracy at 100%, at the offset angle of 60°. We also observe a bit
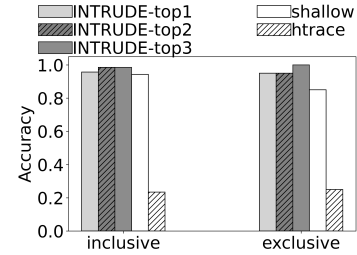


Figure 17: Identification accuracy of INTRUDE when additional participants are used exclusively as victims.

of fluctuation when the offset angle is large. This is because head orientation estimation becomes less stable in these cases. More evaluation data can be used to smooth out the variance.

**Participants Who Are Not in the Training Set.** As mentioned in §6.1, the participants in the above experiments took part in both tasks (training the HMD and playing the role of the victim). To evaluate how the knowledge about the procedures of the former task influences the results of the latter task, we additionally recruited four new participants who exclusively played the victim role and trained the HME only via data from the existing 31 participants. As shown in Figure 17, INTRUDE achieves top-1, 2, 3 accuracies of 95%, 95%, and 100% when the participants are exclusively the victims. The performance is comparable to the inclusive case where the participants appeared in both tasks. This indicates involving participants in both tasks does not affect INTRUDE's performance, validating the appropriateness of our experiment methodology.

## 7.5 Open-World Video Identification

We have evaluated video identification in the in-library scenario, where the video viewed by the victim is known to be from the library. This scenario has been assumed in the majority of prior video identification attacks [31, 56, 69, 95]. However, another scenario exists in practice – open-world video identification, where we do not know if the video being viewed is in the library. The accuracy of open-world identification has been evaluated by theoretically estimating the accuracy drop using a statistical model [75]. We use this approach to evaluate the degradation of INTRUDE in the open-world scenario as the number of out-of-library videos increases.

To model the accuracy drop statistically, the Bayesian Detection Rate (BDR) can be used. BDR estimates the ratio of successful identification among all identification attempts and can be formally expressed as,

$$\text{BDR} = \frac{\text{TPR} \times base}{\text{TPR} \times base + \text{FPR} \times (1 - base)} \quad (4)$$

where TPR (True Positive Rate) is the top-1 accuracy of in-library identification reported above. FPR (False Positive Rate) can be derived from TPR via $FPR = (1 - TPR)\frac{P}{N}$, where $N$ and $P$ are the numbers of negative/positive matches (§6.2) INTRUDE encounters during in-library identification, respectively. The base rate (*base*) is the probability of a test video being in-library, i.e.,

$$base = \frac{\text{Number of in-lib videos}}{\text{Number of both in-lib and out-of-lib videos}} \quad (5)$$

Based on our evaluation results, INTRUDE achieves TPR of 0.96 and FPR of 0.000068. Given *base* = 0.001, the open-world dataset (in-library plus out-of-library videos) is 1,000 times larger than the target library, totaling 635,000 videos. In this case, the BDR is 93%. Given *base* = 0.00025 when the open-world dataset expands to 2,540,000 videos, the BDR achieves 78%. Our analysis indicates INTRUDE still poses great threats even when there are several million 360° videos in the open world.

# 8 Discussion

## 8.1 Limitation and Potential Improvement

**Re-training Cost.** Launching INTRUDE on a new HMD requires the attackers to offline collect recordings of the target HMD and retrain the head orientation estimation model for INTRUDE to recognize the HMD's shape and color. The attackers can capture the recordings of their own head movements when using the target HMD. The brand and model of the HMD can be visually recognized when targeting the victim. Despite the extra model training, the retraining cost is not prohibitive. First, there are only a handful of HMDs on the market. We only need to repeat the training a few times based on demand to make INTRUDE generally applicable. Second, the retraining data can be collected by fewer subjects using the HMD in different outfits. Along with the data augmentation we used (§6.2) and generating additional computer-synthesized data [16], the retraining cost can be minimized.

Notably, the trace-fingerprint matching model of INTRUDE is not HMD-dependent. The current model can be directly used and no retraining is needed when applying INTRUDE to a new HMD. The reason is that the trace-fingerprint matching model takes extracted head movement traces and video fingerprints as input rather than the recording of the new HMD. The TFM can thus infer the video titles in an HMD-agnostic way.

**Non-continuous Video Playback.** We have so far considered continuously-playing videos. This is the most common scenario for 360° video viewing in HMDs because, unlike 2D video viewing with control shortcuts on keyboards, users tend to continuously immerse themselves in 360° video scenes without interruption [63]. However, scenarios where a user pauses, rewinds, or fast-forwards the video can occur, disrupting the continuity of head movements and potentially affecting the attack's success. Fortunately, these infrequent instances can be identified due to its distinctive interactive sequence. For example, pausing necessitates two clicks (one to pause, another to resume). Fast-forwarding or rewinding may require either a single click on the progress bar to skip content or three clicks (pause, forwarding/rewinding button, resume). An alignment module can then be added to INTRUDE to process the recording such that the victim's head movement is aligned with the continuous video playback. In the case of a pause, this new module would discard the recording period while the video is paused. For rewind or fast-forward actions, it extracts the head movement trace once the playback resumes and temporally shifts this trace forward or backward to a video-playing position where the revised trace matches the saliency fingerprints with the highest confidence scores. Similarly, this alignment module can also be used to align the video's start time and the moment when the user clicks the play button if there is a small delay in between.

## 8.2 Mitigations

INTRUDE is effective because video information is always leaking as long as attackers can observe the victim. An approach to prevent the attack is eliminating the direct line of sight between the attacker and the victim [51]. Users could find a large non-transparent obstacle between them and the crowds or turn their heads away from directions that may expose their head movements. However, these efforts are burdensome for casual VR usage and may not always be feasible in certain environments. Moreover, this prevention approach depends on the user's judgment and vigilance of the surrounding people. It is unlikely to work for advanced attackers who disguise themselves or utilize miniature spy cameras [92].

Besides prevention, one could mitigate the success rate of the attack by complicating the attack. For example, in INTRUDE, the roll axis of the VR coordinate system is aligned with the victim's head orientation at the time when the video starts playing. The attacker exploits this fact to extract the victim's head orientation vector from the recording. To undermine the HME, the VR engine can randomly set the direction of the VR roll axis when the video starts. However, this approach cannot completely eliminate the risk, as the victim will eventually follow a certain head movement pattern influenced by the evolving salient content. The attacker may still identify the correct VR roll axis by spatially shifting the saliency maps to positions with the highest confidence scores.

More importantly, users will feel disoriented and have motion sickness every time a video starts playing because the VR coordinate system is not aligned with the head orientation. A full-scale study of the defense capability and usability is needed before this mitigation can be applied.

## 9 Related Work

**Recording-based Side-channel Attacks.** Capturing recordings through the visual channel has been used for recovering keystrokes and unlock patterns on mobile devices such as smartphones and tablets. Existing works have recorded victims (including body movements [21, 43, 46], eyes [20, 82], fingertips [14, 77, 98, 101]), reflective surfaces [11, 12, 66, 96], and backside motion of tablets [82] to perform inference attacks. Recently, there are works showing that sensitive information such as keystrokes entered by participants can be retrieved from virtual backgrounds [35, 73], or reflections from participants' glasses [52, 89].

There is one existing work [47] exploring recording-based side-channel attacks to infer the passwords or keystrokes from users' head movements in HMDs. The attack used optical flow techniques to track the rotation movement of the head in the recording. However, it can only be used for analyzing a very short duration (i.e., during keystrokes) and limited rotation range (i.e., within the virtual keyboard), because this technique has low accuracy and is susceptible to noises. Therefore, they cannot be applied to the case of 360° video viewing, where we must continuously track the head orientation (omnidirectional) over a long time (up to 60s in our case). In this paper, we utilize recordings of head movements of HMD users to launch 360° video identification attacks. Instead of harnessing an unsupervised hand-crafted technique used in [47], we leverage deep learning techniques to train the head orientation estimation model. This enables INTRUDE to achieve state-of-the-art estimation performance while being able to track the user's head movements in 360° video viewing, where the user may look in any direction. To the best of our knowledge, we are the first to use the recording-based visual channel to identify 360° videos in VR HMDs. **Side-channel Attacks on HMDs.** Due to the increasing popularity of VR/AR devices, side-channel attacks on HMDs are gaining more attention. To date, existing side-channel attacks on VR/AR devices [9, 47, 53–55, 58, 78, 93, 103] mainly aim to infer the keystrokes when the user is wearing an HMD. There is one recent work [78] utilizes head motions to infer keystrokes in AR/VR devices, but it relies on a pre-installed malicious app to record sensor readings directly. VR-Spy [9] utilizes channel state information of WiFi signals to recognize keystrokes in VR headsets. Similarly, keystroke inference attacks have been performed on AR devices, e.g., Microsoft Hololens [53, 58]. Meteriz *et al.* [58] explore the AR hand tracking features in

Microsoft Hololens, while Luo *et al.* [53] use the six degree-of-freedom motion tracking information. Unlike these works, the goal of INTRUDE is not to infer the keystrokes, but rather the 360° video being viewed in the HMD.

**Side-channel Attacks for Video Identification.** Among different side-channel techniques to perform video identification attacks, traffic analysis is the most popular type due to the popularity of online video streaming platforms such as Netflix and the ease of accessing the encrypted network traffic in wireless networks. To date, most of the existing works focus on inferring traditional 2D videos [23, 31, 45, 48, 49, 69, 70, 74, 75, 79]. Recently, Bae *et al.* [13] show that video identification attack is feasible in Long Term Evolution (LTE) networks by utilizing broadcast radio signals. Other side channels such as powerline electromagnetic emancipation [26], reflective lights on windows [95], power measurements at the smartphone charging hub [7], and luminance of smart light bulbs [56] have also been utilized to identify 2D videos.

As for 360° videos, a recent study [41] distinguishes 360° videos from regular videos (i.e., a binary classification) for mobile network providers by using packet- and flow-level traffic of encrypted streams. By contrast, INTRUDE is a side-channel attack that identifies 360° video tiles while they are being viewed in HMDs. We harness a novel side channel, i.e., the head movement, which does not exist when viewing conventional 2D videos on laptops and smartphones.

## 10 Conclusion

In this paper, we investigate the leakage of 360° video titles from VR HMD users. In contrast to the public impression that HMDs conceal the displayed content, we show that INTRUDE can infer video titles with high accuracy by recording victims whose vision is blocked and leveraging the unique side channel of head movements. In a typical indoor setup with the camera a few meters away from the victim, INTRUDE can identify video titles with top-1,2,3 accuracy of 96%, 99%, and 99%, respectively. Furthermore, INTRUDE performs well under various noises in different recording distances, lighting, background scenes, and camera placement angles, as well as in the open-world identification scenario with out-of-library videos.

## 11 Acknowledgement

# References

[1] "Airports tap into virtual reality for passenger entertainment," https://www.airport-technology.com/analysis/virtual-reality-passenger-entertainment/.

[2] "How washington's last remaining video rental store changed the course of privacy law," https://www.washingtonpost.com/news/the-switch/.

[3] "Scam emails demand bitcoin, threaten blackmail," https://www.consumer.ftc.gov/blog/2020/04/.

[4] "Tips for launching a public library virtual reality program," https://statetechmagazine.com/article/2019/07/tips-launching-public-library-virtual-reality-program-perfcon.

[5] "U.s. general services administration," http://www.gsa.gov/portal/content/101308.

[6] A. D. Abreu, C. Ozcinar, and A. Smolic, "Look around you: Saliency maps for omnidirectional images in vr applications," in *IEEE International Conference on Quality of Multimedia Experience (QoMEX)*, May 2017.

[7] S. Acharya, A. Serwadda, and A. V. Bilbao, "That phone charging hub knows your video playlist!" in *2021 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computing, Scalable Computing & Communications, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/IOP/SCI)*. IEEE, 2021, pp. 160–169.

[8] S. Afzal, J. Chen, and K. Ramakrishnan, "Characterization of 360-degree videos," in *Proceedings of the Workshop on Virtual Reality and Augmented Reality Network*, 2017, pp. 1–6.

[9] A. Al Arafat, Z. Guo, and A. Awad, "Vr-spy: A side-channel attack on virtual key-logging in vr headsets," in *2021 IEEE Virtual Reality and 3D User Interfaces (VR)*. IEEE, 2021, pp. 564–572.

[10] L. Argyriou, D. Economou, and V. Bouki, "Design methodology for 360 immersive video applications: the case study of a cultural heritage virtual tour," *Personal and Ubiquitous Computing*, vol. 24, pp. 843–859, 2020.

[11] M. Backes, T. Chen, M. Dürmuth, H. P. Lensch, and M. Welk, "Tempest in a teapot: Compromising reflections revisited," in *2009 30th IEEE Symposium on Security and Privacy*. IEEE, 2009, pp. 315–327.

[12] M. Backes, M. Dürmuth, and D. Unruh, "Compromising reflections-or-how to read lcd monitors around the corner," in *2008 IEEE Symposium on Security and Privacy (sp 2008)*. IEEE, 2008, pp. 158–169.

[13] S. Bae, M. Son, D. Kim, C. Park, J. Lee, S. Son, and Y. Kim, "Watching the watchers: Practical video identification attack in lte networks," in *USENIX Security Symposium*, 2022.

[14] D. Balzarotti, M. Cova, and G. Vigna, "Clearshot: Eavesdropping on keyboard input from video," in *2008 IEEE Symposium on Security and Privacy (sp 2008)*. IEEE, 2008, pp. 170–183.

[15] M. Ben-Ari, "A tutorial on euler angles and quaternions," *Weizmann Institute of Science, Israel*, 2014.

[16] C. Bermejo, D. Chatzopoulos, and P. Hui, "Eyeshopper: Estimating shoppers' gaze using cctv cameras," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2765–2774.

[17] N. Carlsson and D. Eager, "Cross-user similarities in viewing behavior for 360 video and caching implications," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 2021.

[18] C. Chen, Q. Chen, J. Xu, and V. Koltun, "Learning to see in the dark," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3291–3300.

[19] T. Chen, L. Xu, X. Xu, and K. Zhu, "Gestonhmd: Enabling gesture-based interaction on low-cost vr head-mounted display," *IEEE Transactions on Visualization and Computer Graphics*, vol. 27, no. 5, pp. 2597–2607, 2021.

[20] Y. Chen, T. Li, R. Zhang, Y. Zhang, and T. Hedgpeth, "Eyetell: Video-assisted touchscreen keystroke inference from eye movements," in *2018 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2018, pp. 144–160.

[21] Y. Chen, Y. Du, C. Xu, Y. Yu, H. Liu, H. Dai, Y. Ren, and J. Yu, "Armspy: Video-assisted pin inference leveraging keystroke-induced arm posture changes," in *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*. IEEE, 2022, pp. 1878–1887.

[22] X. Corbillon, F. De Simone, and G. Simon, "360-degree video head movement dataset," in *Proceedings of the 8th ACM on Multimedia Systems Conference (MMSys'17)*, 2017.

[23] R. Dubin, A. Dvir, O. Pele, and O. Hadar, "I know what you saw last minute—encrypted http adaptive video streaming title classification," *IEEE transactions on information forensics and security*, vol. 12, no. 12, pp. 3039–3049, 2017.

[24] P. e Pubs, "Is your library ready for the reality of virtual reality?" https://docs.lib.purdue.edu/cgi/viewcontent.cgi?article=2081&context=charleston.

[25] EducationWeek, "Vr devices collect 'intimate' data, lack privacy protections. should schools invest?" https://www.edweek.org/technology/vr-devices-collect-intimate-data-lack-privacy-protections-should-schools-invest/2022/11.

[26] M. Enev, S. Gupta, T. Kohno, and S. N. Patel, "Televisions, video privacy, and powerline electromagnetic interference," in *Proceedings of the 18th ACM conference on Computer and communications security*, 2011, pp. 537–550.

[27] C. Fan, J. Lee, W. Lo, C. Huang, K. Chen, and C. Hsu, "Fixation prediction for 360 video streaming in head-mounted virtual reality," in *ACM Workshop on Network and Operating Systems Support for Digital Audio and Video (NOSSDAV)*, Jun. 2017.

[28] T. Feigl, C. Mutschler, and M. Philippsen, "Supervised learning for yaw orientation estimation," in *2018 international conference on indoor positioning and indoor navigation (IPIN)*. IEEE, 2018, pp. 206–212.

[29] E. M. Foxlin, M. Harrington, and Y. Altshuler, "Miniature six-dof inertial system for tracking hmds," in *Helmet-and Head-Mounted Displays III*, vol. 3362. International Society for Optics and Photonics, 1998, pp. 214–228.

[30] K. Grinyer and R. J. Teather, "Effects of field of view on dynamic out-of-view target search in virtual reality," in *2022 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. IEEE, 2022, pp. 139–148.

[31] J. Gu, J. Wang, Z. Yu, and K. Shen, "Traffic-based side-channel attack in video streaming," *IEEE/ACM Transactions on Networking*, vol. 27, no. 3, pp. 972–985, 2019.

[32] A. J. Hanson, "Visualizing quaternions," pp. 1–es, 2005.

[33] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.

[34] ——, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[35] J. M. Hilgefort, D. Arp, and K. Rieck, "Spying through virtual backgrounds of video calls," in *Proceedings of the 14th ACM Workshop on Artificial Intelligence and Security*, 2021, pp. 135–144.

[36] M.-C. Hsieh and J.-J. Lee, "Preliminary study of vr and ar applications in medical and healthcare education," *J Nurs Health Stud*, vol. 3, no. 1, p. 1, 2018.

[37] L. Jiang, M. Xu, X. Wang, and L. Sigal, "Saliency-guided image translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 509–16 518.

[38] M. Jiang, S. Huang, J. Duan, and Q. Zhao, "Salicon: Saliency in context," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1072–1080.

[39] Kaspersky, "What are the security and privacy risks of vr and ar," https://usa.kaspersky.com/resource-center/threats/security-and-privacy-risks-of-ar-and-vr.

[40] F. Katsuki and C. Constantinidis, "Bottom-up and top-down attention: different processes and overlapping neural systems," *The Neuroscientist*, vol. 20, no. 5, pp. 509–521, 2014.

[41] C. Kattadige, A. Raman, K. Thilakarathna, A. Lutu, and D. Perino, "360norvic: 360-degree video classification from mobile encrypted video traffic," in *Proceedings of the 31st ACM Workshop on Network and Operating Systems Support for Digital Audio and Video*, 2021, pp. 58–65.

[42] P. Kellnhofer, A. Recasens, S. Stent, W. Matusik, and A. Torralba, "Gaze360: Physically unconstrained gaze estimation in the wild," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6912–6921.

[43] H. Khan, U. Hengartner, and D. Vogel, "Evaluating attack and defense strategies for smartphone pin shoulder surfing," in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, 2018, pp. 1–10.

[44] L. Leveque and H. Liu, "An eye-tracking database of video advertising," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 425–429.

[45] Y. Li, Y. Huang, R. Xu, S. Seneviratne, K. Thilakarathna, A. Cheng, D. Webb, and G. Jourjon, "Deep content: Unveiling video streaming content from encrypted wifi traffic," in *2018 IEEE 17th International Symposium on Network Computing and Applications (NCA)*. IEEE, 2018, pp. 1–8.

[46] Z. Li, Q. Sun, Y. Lian, and D. D. Giusto, "An association-based graphical password design resistant to shoulder-surfing attack," in *2005 IEEE international conference on multimedia and expo*. IEEE, 2005, pp. 245–248.

[47] Z. Ling, Z. Li, C. Chen, J. Luo, W. Yu, and X. Fu, "I know what you enter on gear vr," in *2019 IEEE Conference on Communications and Network Security (CNS)*. IEEE, 2019, pp. 241–249.

[48] Y. Liu, C. Ou, Z. Li, C. Corbett, B. Mukherjee, and D. Ghosal, "Wavelet-based traffic analysis for identifying video streams over broadband networks," in *IEEE GLOBECOM 2008-2008 IEEE Global Telecommunications Conference*. IEEE, 2008, pp. 1–6.

[49] Y. Liu, A.-R. Sadeghi, D. Ghosal, and B. Mukherjee, "Video streaming forensic–content identification with traffic snooping," in *International Conference on Information Security*. Springer, 2010, pp. 129–135.

[50] W.-C. Lo, C.-L. Fan, and J. Lee, "360-degree video viewing dataset in head-mounted virtual reality," in *Proceedings of the 8th ACM on Multimedia Systems Conference (MMSys'17)*, 2017.

[51] J. Long, *No tech hacking: A guide to social engineering, dumpster diving, and shoulder surfing*. Syngress, 2011.

[52] Y. Long, C. Yan, S. Prasad, W. Xu, and K. Fu, "Private eye: On the limits of textual screen peeking via eyeglass reflections in video conferencing," *arXiv preprint arXiv:2205.03971*, 2022.

[53] S. Luo, X. Hu, and Z. Yan, "Holologger: Keystroke inference on mixed reality head mounted displays," in *2022 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. IEEE, 2022, pp. 445–454.

[54] S. Luo, A. Nguyen, H. Farooq, K. Sun, and Z. Yan, "Eavesdropping on controller acoustic emanation for keystroke inference attack in virtual reality," in *The Network and Distributed System Security Symposium (NDSS)*, Feb. 2024.

[55] S. Luo, A. Nguyen, C. Song, F. Lin, W. Xu, and Z. Yan, "Oculock: Exploring human visual system for authentication in virtual reality head-mounted display," in *2020 The Network and Distributed System Security Symposium (NDSS)*, 2020.

[56] A. Maiti and M. Jadliwala, "Light ears: Information leakage via smart lights," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 3, no. 3, pp. 1–27, 2019.

[57] R. McGloin, K. Farrar, and M. Krcmar, "Video games, immersion, and cognitive aggression: does the controller matter?" *Media psychology*, vol. 16, no. 1, pp. 65–87, 2013.

[58] Ü. Meteriz-Yıldıran, N. F. Yıldıran, A. Awad, and D. Mohaisen, "A keylogging inference attack on air-tapping keyboards in virtual environments," in *2022 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. IEEE, 2022, pp. 765–774.

[59] E. Murphy-Chutorian and M. M. Trivedi, "Head pose estimation in computer vision: A survey," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 4, pp. 607–626, 2008.

[60] A. Nassani, L. Zhang, H. Bai, and M. Billinghurst, "Showmearound: Giving virtual tours using live 360 video,"

in *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021, pp. 1–4.

[61] A. Nguyen and Z. Yan, "A saliency dataset for 360-degree videos," in *Proceedings of the 10th ACM Multimedia Systems Conference*, 2019, pp. 279–284.

[62] A. Nguyen, Z. Yan, and K. Nahrstedt, "Your attention is unique: Detecting 360-degree video saliency in head-mounted display for head movement prediction," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 1190–1198.

[63] C. Oh, F. Herrera, and J. Bailenson, "The effects of immersion and real-world distractions on virtual social interactions," *Cyberpsychology, Behavior, and Social Networking*, vol. 22, no. 6, pp. 365–372, 2019.

[64] M. Patacchiola and A. Cangelosi, "Head pose estimation in the wild using convolutional neural networks and adaptive gradient methods," *Pattern Recognition*, vol. 71, pp. 132–143, 2017.

[65] J. Radianti, T. A. Majchrzak, J. Fromm, and I. Wohlgenannt, "A systematic review of immersive virtual reality applications for higher education: Design elements, lessons learned, and research agenda," *Computers & Education*, vol. 147, p. 103778, 2020.

[66] R. Raguram, A. M. White, D. Goswami, F. Monrose, and J.-M. Frahm, "ispy: automatic reconstruction of typed input from compromising reflections," in *Proceedings of the 18th ACM conference on Computer and communications security*, 2011, pp. 527–536.

[67] B. Ray, J. Jung, and M.-C. Larabi, "A low-complexity video encoder for equirectangular projected 360 video content," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 1723–1727.

[68] Redditor-crazy_goat, "Virtual reality security and privacy," https://www.reddit.com/r/oculus/comments/8k42kp/my_experience_using_oculus_go_on_a_35_hour_flight/n.

[69] A. Reed and B. Klimkowski, "Leaky streams: Identifying variable bitrate dash videos streamed over encrypted 802.11 n connections," in *2016 13th IEEE Annual Consumer Communications & Networking Conference (CCNC)*. IEEE, 2016, pp. 1107–1112.

[70] A. Reed and M. Kranch, "Identifying https-protected netflix videos in real-time," in *Proceedings of the Seventh ACM on Conference on Data and Application Security and Privacy*, 2017, pp. 361–368.

[71] RoadToVR, "Meta has sold nearly 20 million quest headsets, but retention struggles remain," https://www.roadtovr.com/quest-sales-20-million-retention-struggles/.

[72] N. Ruiz, E. Chong, and J. M. Rehg, "Fine-grained head pose estimation without keypoints," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 2074–2083.

[73] M. Sabra, A. Maiti, and M. Jadliwala, "Background buster: Peeking through virtual backgrounds in online video calls," 2022.

[74] T. S. Saponas, J. Lester, C. Hartung, S. Agarwal, T. Kohno *et al.*, "Devices that tell on you: Privacy trends in consumer ubiquitous computing." in *USENIX Security Symposium*, 2007, pp. 55–70.

[75] R. Schuster, V. Shmatikov, and E. Tromer, "Beauty and the burst: Remote identification of encrypted video streams," in *26th {USENIX} Security Symposium ({USENIX} Security 17)*, 2017, pp. 1357–1374.

[76] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of big data*, vol. 6, no. 1, pp. 1–48, 2019.

[77] D. Shukla, R. Kumar, A. Serwadda, and V. V. Phoha, "Beware, your hands reveal your secrets!" in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, 2014, pp. 904–917.

[78] C. Slocum, Y. Zhang, N. Abu-Ghazaleh, and J. Chen, "Going through the motions: Ar/vr keylogging from user head motions," in *USENIX Security*, 2023.

[79] J. Song, S. Lee, B. Kim, S. Seol, B. Lee, and M. Kim, "Vtim: Video title identification using open metadata," *IEEE Access*, vol. 8, pp. 113 567–113 584, 2020.

[80] M. Speicher, K. Lewis, and M. Nebeling, "Designers, the stage is yours! medium-fidelity prototyping of augmented & virtual reality interfaces with 360theater," *Proceedings of the ACM on Human-Computer Interaction*, vol. 5, no. EICS, pp. 1–25, 2021.

[81] Statista, "Virtual reality (vr) - statistics & facts," https://www.statista.com/topics/2532/virtual-reality-vr/.

[82] J. Sun, X. Jin, Y. Chen, J. Zhang, Y. Zhang, and R. Zhang, "Visible: Video-assisted keystroke inference from tablet backside motion." in *NDSS*, 2016.

[83] T. Verge, "Surveillance will follow us into the 'metaverse', and our bodies could be its new data source," https://www.washingtonpost.com/technology/2022/01/13/privacy-vr-metaverse/.

[84] ——, "This is meta's ar / vr hardware roadmap for the next four years," https://www.theverge.com/2023/2/28/23619730/meta-vr-oculus-ar-glasses-smartwatch-plans.

[85] Vittorio, "In-flight vr: Using the oculus quest on a plane," https://medium.com/@vibronet/in-flight-vr-using-the-oculus-quest-on-a-plane-38c9808c32b2.

[86] J.-N. Voigt-Antons, T. Kojic, D. Ali, and S. Möller, "Influence of hand tracking as a way of interaction in virtual reality on user experience," in *2020 Twelfth International Conference on Quality of Multimedia Experience (QoMEX)*. IEEE, 2020, pp. 1–4.

[87] E. Von Zezschwitz, A. De Luca, P. Janssen, and H. Hussmann, "Easy to draw, but hard to trace? on the observability of grid-based (un) lock patterns," in *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, 2015, pp. 2339–2342.

[88] W. Wang, J. Shen, and F. Porikli, "Saliency-aware geodesic video object segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3395–3402.

[89] H. Wasswa and A. Serwadda, "The proof is in the glare: On the privacy risk posed by eyeglasses in video calls," in *Proceedings of the 2022 ACM on International Workshop on Security and Privacy Analytics*, 2022, pp. 46–54.

[90] S. Weisberg, *Applied linear regression*. John Wiley & Sons, 2005, vol. 528.

[91] C. Wu, Z. Tan, and Z. Wang, "360-degree video viewing dataset in head-mounted virtual reality," in *Proceedings of the 8th ACM on Multimedia Systems Conference (MMSys'17)*, 2017.

[92] L. Wu, X. Du, and X. Fu, "Security threats to mobile multimedia applications: Camera-based attacks on mobile phones," *IEEE Communications Magazine*, vol. 52, no. 3, pp. 80–87, 2014.

[93] Y. Wu, C. Shi, T. Zhang, P. Walker, J. Liu, N. Saxena, and Y. Chen, "Privacy leakage via unrestricted motion-position sensors in the age of virtual reality: A study of snooping typed input on virtual keyboards," in *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 2023, pp. 3382–3398.

[94] XRToday, "Virtual reality security and privacy," https://www.xrtoday.com/virtual-reality/virtual-reality-security-and-privacy/.

[95] Y. Xu, J.-M. Frahm, and F. Monrose, "Watching the watchers: Automatically inferring tv content from outdoor light effusions," in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, 2014, pp. 418–428.

[96] Y. Xu, J. Heinly, A. M. White, F. Monrose, and J.-M. Frahm, "Seeing double: Reconstructing obscured typed input from repeated compromising reflections," in *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, 2013, pp. 1063–1074.

[97] G. Ye, Z. Tang, D. Fang, X. Chen, K. I. Kim, B. Taylor, and Z. Wang, "Cracking android pattern lock in five attempts," in *Proceedings of the 2017 Network and Distributed System Security Symposium 2017 (NDSS 17)*. Internet Society, 2017.

[98] G. Ye, Z. Tang, D. Fang, X. Chen, W. Wolff, A. J. Aviv, and Z. Wang, "A video-based attack for android pattern lock," *ACM Transactions on Privacy and Security (TOPS)*, vol. 21, no. 4, pp. 1–31, 2018.

[99] Youtube, "Youtube features - virtual reality," https://www.youtube.com/@360/feature.

[100] T. Yubing, F. A. Cheikh, F. F. E. Guraya, H. Konik, and A. Trémeau, "A spatiotemporal saliency model for video surveillance," *Cognitive Computation*, vol. 3, pp. 241–263, 2011.

[101] Q. Yue, Z. Ling, X. Fu, B. Liu, K. Ren, and W. Zhao, "Blind recognition of touched keys on mobile devices," in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, 2014, pp. 1403–1414.

[102] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.

[103] Y. Zhang, C. Slocum, J. Chen, and N. Abu-Ghazaleh, "It's all in your head (set): Side-channel attacks on ar/vr systems," in *USENIX Security*, 2023.

[104] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li, "Face alignment across large poses: A 3d solution," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 146–155.

# A  Appendix

To accurately track the head movements of victims in camera recordings, it is crucial for the Head Movement Estimation (HME) scheme to precisely estimate head orientations at a granular level. This requires the HME training set to include all possible head orientations that victims may exhibit. In this section, we explain the comprehensive collection process for all common head orientation angles for HME training. We begin with a basic explanation of the VR coordinate system, necessary for understanding the data labeling procedure, followed by a detailed description of the experimental setups for collecting training set samples. Finally, we outline the annotation process required to generate the ground truth.

## A.1  Specifications of Coordinates in Virtual Reality

The VR OS tracks user head movements using the VR coordinate system, a critical factor in rendering objects in VR. Every head movement triggers the system to adjust the positions of virtual objects relative to the viewport, thereby creating the illusion of user exploration within the virtual environment. The system represents head orientation as a rotational operation. This operation rotates a unit vector, i.e., a vector where one component equals 1 while the remaining components equal zero, from the reference axis (the roll axis in Figure 7) to the current orientation of the head. Many VR engines use quaternions, a tuple of 4 elements equivalent to a $3 \times 3$ rotation matrix, to represent the rotation operation. The quaternion is preferred as the rotation matrix due to its compactness [32]. As such, the 3D head orientation vector can be derived simply by applying the rotational operation to the unit vector from the reference axis.

## A.2  Experiment Setup

This section discusses the setup and procedure to collect samples for the training set. The data collection procedure includes three main steps, designing a virtual room, designing a viewing procedure, and capturing camera recordings.

### A.2.1  Design of the Virtual Room

We design a virtual room that enables subjects to exhibit all aspects of their head and HMD movements to the recording

Figure 18: Design of the virtual room. a) the positions of cubes around a subject b) a view inside the virtual room.

camera (*RC*). The virtual room consists of five cubes placed around the subject, including two colored cubes and three gray cubes. The cubes' arrangement is depicted in Figure 18. All cubes are initially suspended 4.5 meters above the floor and descend at a constant speed. After 30 seconds, the cubes reach the floor and subsequently disappear. Figure 18b shows part of the virtual room with 3 gray cubes visible in the viewport.

### A.2.2 Viewing Procedure

After putting on the HMD, subjects are asked to look around in a systematic way such that *RC* could capture their heads at different angles. In particular, subjects turn their heads from the red cube, glance over the gray cubes, and stop briefly at the blue cube before starting over from the blue cube to the red cube. As the head follows the dropping cubes, the *RC* placed at a distance away captures their head poses. After 30 seconds, the cube drops to the floor, subjects are told to look around and explore locations in the room that have not been seen in the first 30 seconds. This procedure lasts for 60 seconds.

### A.2.3 Capturing the Camera Recordings

We used the same set of subjects and followed the same setup described in §6 to capture camera recordings. The virtual room was installed in the HMD. Three *RCs* were positioned in front of the participants at distances of 1.5m, 3.0m, and 4.5m. Participants were instructed to wear the HMD and follow the pre-established viewing procedure while the *RC* captured the camera recordings. Simultaneously, the timestamped head movement data was collected from the HMD gyroscope sensors and stored in head movement logs. These logs will be used for subsequent annotation. This procedure resulted in camera recordings capturing head poses of subjects ranging from $-179°$ to $179°$ in yaw and from $-55°$ to $55°$ in pitch.

### A.3 Labeling the Camera Recordings

The camera recordings were divided into individual frames, with each frame's head orientation annotated. Following this,

we identified and removed a source of noise from the annotated data.

### A.3.1 Data Annotation

Our method for annotating head poses in each frame of the camera recordings is both economical and precise, thanks to the availability of the gyroscope sensors within the HMD. To generate a large volume of annotated data, the recordings and the gyroscope data from the head movement logs were synchronized. Upon synchronization, every frame was correctly linked to a head orientation vector from the head movement log. The primary challenge was determining the time offset between the camera recording and the head movement log. Correctly calculating this offset requires the annotator to manually identify a key frame $f_i$ at time $t_i$ from the camera recording and the associated key head orientation $v_j$ at time $t_j$ in the head movement log. Let $((f_i, t_i), (v_j, t_j))$ denote the selected pair before the synchronization. Then, the offset is the time interval between the two key points, $t_i - t_j$. After being synchronized, $(f_i, v_j, t_j)$ is the annotated frame with the ground truth head orientation vector and the adjusted timestamp. The time taken to process this step for each video lasts only a few minutes.

### A.3.2 Noise Removal

In this section, we identify yaw drift bias as a source of noise in the annotated head orientation and discuss our strategy for its elimination.

**Noise source.** Gyroscope data was utilized to annotate head orientations for each frame. The gyroscope measures the angular rotational forces of the head in three axes (yaw, pitch, and roll), enabling the virtual coordinates to remain static regardless of where the user looks. This is vital for accurate head orientation tracking over time. However, the yaw component of the gyroscope data is prone to noise. Consequently, the VR coordinate system slowly drifts along the yaw axis. This form of noise is referred to as yaw drift bias [29]. The yaw drift error accumulates over the course of viewing sessions, resulting in inconsistencies in the labeled data from later parts of the camera *recordings*. As such, yaw drift bias significantly affects the accuracy of the annotated head poses in each camera recording.

The yaw drift can be detected by observing the inconsistency among two video frames at two different times. Assuming $f_i, f_{i'}$ are two frames at time $t_i$ and $t_{i'}$, and $v_i, v_{i'}$ are their annotated head orientation vectors recorded from the gyroscope, respectively. Assuming the head orientations captured by $f_i$ and $f_{i'}$ are similar, e.g., the victim looks directly at the camera. If the yaw drift occurs, there is an inconsistency between the annotated frame, i.e., $f_i - f_{i'} \neq 0$.

**Noise removal strategy.** One method to counteract this bias involves modeling the yaw drift as a linear function of time

and subtracting the yaw drift component from the annotated data. Linear modeling has been widely used for approximating data in regression and forecasting tasks [90]. In our case, it is possible to approximate the function parameters from two data points, $(f_i, v_i)$ and $(f_{i'}, v_{i'})$. Let $YD = \theta x + \theta_0$ be the linear function modeling the yaw drift. If the component of a 3D head orientation vector $v_i$ is represented as $(v_{i,1}, v_{i,2}, v_{i,3})$, then the parameter $\theta$ and $\theta_0$ can be calculated from the head orientation vectors $v_i$ and $v_{i'}$ as:

$$\theta = (v_{i',1} - v_{i,1})/(t_{i'} - t_i)$$

$$\theta_0 = v_{i,1}$$

Eliminating the bias involves two fundamental quaternion operations, $create\_quat(axis, deg)$ and $rotate(q, v)$ [15]. $create\_quat(axis, deg)$ is the function to build a rotation operation rotating a vector for $deg$ degrees around an axis specified by the unit vector $axis$. $rotate(q, v)$ is the function returning a 3D vector after applying the rotation operation $q$ on the vector $v$. Let $[v_1, .., v_k, .., v_N]$ be the annotated head orientations of a camera recording. For every head orientation vector $v_k$, a rotation operation $q_i = create\_quat([1, 0, 0], -YD(t_i))$ is constructed to rotate $v_k$ for $-YD(t_i)$ degrees along the yaw axis. The resulting $u_k = rotate(q_i, v_i)$ is the head orientation vector after the yaw drift bias is eliminated.