

Video-side Channel Attack for Android Pattern Lock

Guixin Ye, Zhanyong Tang, Dingyi Fang, Xiaojiang Chen, Willy Wolf and Zheng Wang

Abstract—Pattern lock is widely used as a mechanism for authentication and authorization on Android devices. This paper presents a novel video-based attack to reconstruct Android lock patterns from video footage filmed using a mobile phone camera. Unlike prior attacks on pattern lock, our approach does not require the video to capture ant content displayed on the screen. Instead, we employ a computer vision algorithm to track the fingertip movements to infer the pattern. Using the geometry information extracted from the tracked fingertip motions, our approach is able to accurately identify a small number of (often one) candidate patterns to be tested by an adversary. We thoroughly evaluated our approach using 120 unique patterns collected from 215 independent users, by applying it to reconstruct patterns from video footage filmed using mobile phone cameras. Experimental results show that our approach can break over 95% of the patterns in five attempts before the device is automatically locked by the Android operating system. We discovered that, in contrast to many people’s belief, complex patterns do not offer stronger protection under our attacking scenarios. This is demonstrated by the fact that we are able to break all but one complex patterns as opposed to 60% of the simple patterns in the first attempt. Since our threat model is common in day-to-day life, this paper calls for the community to revisit the risks of using Android pattern lock to protect sensitive information.

Index Terms—Pattern lock, Fingertip movement, Video-based attack, Sensitive information.

I. INTRODUCTION

Pattern lock is widely used on Android devices to protect sensitive information. It is preferred by some users over PIN- or text-based passwords, as psychology studies show that the human brain remembers and recalls visual information better than numbers and letters [?]. According to a recent study, 40% of the Android users use patterns to protect their devices instead of a PIN [?]. Pattern lock is also used for authentication. For example, Alipay, the largest third-party online-payment platform, uses pattern lock for login authentication. Given its pervasive usage, a security breach of the pattern lock could lead to serious consequences.

Researchers have demonstrated a number of ways to crack Android pattern lock. For example, smudge attacks use the oily residues left on the screen to recover the pattern [?]. However, this approach relies on the persistence of the smudge which can be easily destroyed by subsequent on-screen activities after unlocking. In a recent study, Zhang *et al.* [?] shows that it is possible to infer a locking pattern by analyzing how the WiFi signal is affected by the finger motions when drawing the pattern. Their approach is restricted to a limit set of scenarios due two reasons: (1) the complex setup of the

attacking environment and (2) the WiFi signal can be disrupted by any moving objects nearby or body movements.

Recently, video-based side-channel attacks are shown to be effective in reconstructing PIN- or text-based passwords. Some of the early work in this area rely on video footage filmed using a camera directly facing the screen or the keyboard [? ?]. Latest work shows that this limit can be lifted by exploiting spatial-temporal dynamics of the hands during typing [?]. Despite the success of video-based attacks on PIN- and text-based passwords, no work so far has exploited video-based side-channels to crack pattern lock. To do so, the attack must address a number of new challenges. These include: How to map the user’s fingertip movements to a graphical structure consisting of continuous points instead of discrete keystrokes? How to transform the fingertip movements tracked from the camera’s perspective to the user’s view point to correctly reconstruct the pattern? How to cancel the camera shake effect that can significantly affect the performance of the attack? how to identify two overlapping line segments of a pattern? The size of the touch-screen or the pattern grid can vary from one device or one application to the other, how can the algorithm adapt to these changes? These new challenges make prior work video-based attacks inapplicable. To overcome these challenges requires creative solutions to be constructed in the new application context of pattern lock.

This paper presents a novel approach to crack Android pattern lock using video footage that captures the user’s fingertip motions when drawing the pattern. Unlike smudge attacks [?], our approach does not require the video footage or images to be captured by a camera directly faced the screen. Furthermore, the video can be filmed at a distance of 2 meters from the user in public places. Such a distance is less likely to raise suspicion compared to shoulder surfing [?] that requires a closer observation distance to have a clear sight of the content displayed on the screen.

Our attack employs a computer vision algorithm to track the fingertip motions from the video. Using the geometry information extracted from the fingertip motions, it then maps the tracked fingertip locations to a small number of (often just one) candidate patterns to be tested on the target device.

We thoroughly evaluate our approach using 120 unique patterns collected from independent users. We show that our approach is effective in inferring candidate patterns and as a result, an attacker can unlock the target device with a success rate of over 95% (up to 97.5%) in five attempts. We demonstrate that, in contrast to many people’s belief, complex patterns do not provide stronger protection over simple patterns under our attack. According to a recent study [?], people tend

G. Ye, Z. Tang, D. Fang, and X. Chen are with Northwest University, China
Willy Wolf and Zheng Wang are with Lancaster University, U.K.

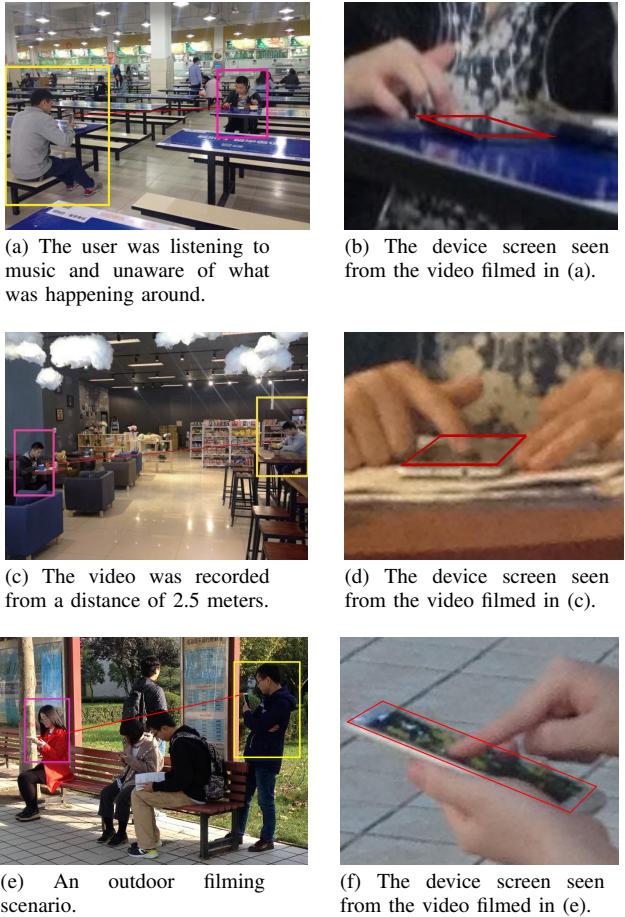


Figure 1. Examples of scenarios in which a mobile phone camera is used to film the unlocking process. In these scenarios, the camera does not need to have a clear sight of the screen.

to use complex patterns for critical applications such as online banking and shopping applications. Our finding suggests that using pattern lock to protect sensitive information could be risky. **FIXED:** Intuitively, our approach is able to crack PIN-based passwords by analyzing the fingertip movement when typing the passwords. We evaluate our approach using 30 four-digital PIN-based passwords and the experimental result shows we can break most of passwords within five attempts.

FIXED: Given the wide usage of pattern lock, we propose a countermeasure to defend against video-based attack without external hardware. This method can effectively confuse the attack to infer correct pattern by disabling a necessary attack factor—user can draw any pattern before or after drawing the correct pattern.

Contributions This work is the first attack to reconstruct graphical patterns from video footage without requiring the video to capture any contents displayed on the screen. Our attack exploits techniques developed in the computer vision domain to address the key challenges highlighted above. The key contribution of this paper is a new attack for Android pattern lock that has not seen in prior work. This paper makes the following specific contributions:

- **A New Attack.** This is the first work to reconstruct locking patterns without relying on the content shown on the screen (Section II-B). Experimental results show that

our method can break over 95% of the lock patterns in five attempts (Section VI-A). Given that the Android operating system (OS) allows five tries before locking the device, our attack represents a real threat for pattern lock.

- **Identifying New Vulnerabilities.** According to a recent study [?], direct observation techniques, e.g. shoulder surfing, are considered to be a low risk due to the close distance between the attacker and the user (in order to gain a clear sight of the device screen). As a result, many users may underestimate the dangers from using pattern lock in public places. Under our attack, filming can be carried out at a distance of 2 meters from the user and the camera does not need to directly face the target device. Such a camera setting makes our attack less likely to raise suspicion and more likely to success when compared to direct observation techniques. For instance, the video can be filmed by an adversary who pretends to interact with his phone, sitting next to the user in a public place (see Figure 1). In many similar scenarios, users will not be suspicious of the attacker’s behavior.

- **New Findings.** Our study suggests that complex patterns are more vulnerable under video-based attacks (Section VI-A). This finding debunks many people’s conception that more complex patterns give stronger protection. Therefore, our work sheds new insights on the practical use of pattern lock.

II. BACKGROUND

A. Android Pattern Lock

Pattern lock is widely used to protect sensitive information and perform authentication on Android touch-screen devices. The motivation for user to prefer to pattern lock is that graphical passwords are easily to be recalled comparing to alphanumeric characters [? ?]. To unlock a device protected with pattern lock, the user is asked to draw a predefined sequence of connected dots on a pattern grid¹. Figure 2 (e) shows a pattern which consists of seven dots on a 3×3 grid. To form a pattern, the user starts by selecting one dot as the starting point and then swiping over multiple dots of the grid until the fingertip is lifted from the screen. There are several rules for creating an Android pattern: (1) a pattern must consist of at least four dots; (2) each dot can only be visited once; and (3) a previously unvisited dot will become visited if it is part of a horizontal, vertical or diagonal line segment of the pattern. Taking into account these constraints, the total number of possible patterns on a 3×3 grid is 389,112 [?]. Given the large number of possible patterns, performing brute-force attacks on Android pattern lock is ineffective, because by default the device will be automatically locked after five failed tries. Further, the patterns with complex structures are hard to guess compared to text-based passwords [? ?].

B. Threat Model

In our threat model, we assume an adversary wants to access some sensitive information from or to install malware on a

¹In this paper we use the Android default pattern grid with 3×3 dots, unless otherwise stated.

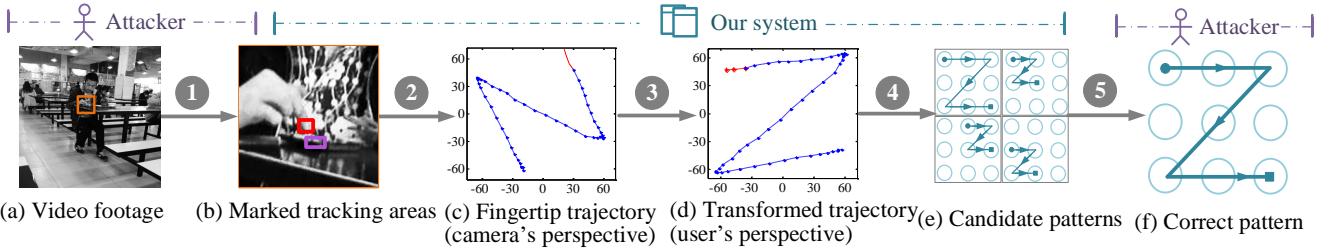


Figure 2. Overview of the attack. Our system takes in a video segment that records the unlocking process (a). The adversary first marks two areas of interest on the first video frame (b): one contains the fingertip involved in pattern drawing, and the other contains part of the device. Our system then tries to track the fingertip’s location w.r.t. to the device. The tracking algorithm produces a fingertip movement trajectory from the camera’s perspective (c) which is then transformed to the user’s perspective (d). Finally, the resulted trajectory in (d) is mapped to several candidate patterns (e) to be tested on the target device (f).

target device that is protected by pattern lock. This type of attacks is mostly likely to be performed by an attacker who can physically access to the target device for a short period of time (e.g. via attending a meeting or party where the user presents). To quickly gain access to the device without raising suspicion, the attacker would like to obtain the user’s locking pattern in advance. The attack starts from filming how the user unlocks the device. Video recording can be done on-site or ahead of time (probably with the help of someone hired by the attack). The video will then be processed to identify a small number of patterns to be tested on the target device. Because filming can be carried out in from a distance of as far as 2.5 meters using a mobile phone camera and the camera does not need to directly face the target device, this activity often will not be noticed by the user. Moreover, given that many users use the same pattern across devices and applications, the pattern obtained from one device could also be used to break the user’s other devices. *The goal of this paper is to demonstrate the feasibility of a new attack and its implications to the use of pattern lock.*

Examples of Filming Scenarios Figure 1 illustrates three scenarios where filming can be performed without raising suspicion to many users. For all the examples presented in Figure 1, the filming camera had a left- or right-front view angle from the target device and did not directly face the screen of the target device. Due to the filming distance (2-3 meters), the recorded video typically does not have a clear vision of the content displayed on the screen. This observation can be confirmed by the video snapshot placing alongside each scenario, where it is impossible to identify the content shown on the screen. The examples given in Figure 1 are some of the day-to-day scenarios where security of the user’s device can be compromised under our attack.

Assumptions Our attack requires the video footage to have a vision of the user’s fingertip involved in pattern drawing and part of the device (e.g. an edge of a phone). We believe this is a reasonable assumption because in practice many users often do not fully cover their fingers and the entire device when drawing a pattern. This is particularly true when holding a large-screen device by hands. To launch the attack, the attacker needs to know the layout of the grid, e.g. whether it is a 3×3 or a 6×6 grid. Our approach is to generate a set of candidate patterns for each of the Android pattern grids and the attacker can simply decide which set of candidate patterns to use after seeing the target device (at the time the layout of the grid will become available). However, unlike prior work on video-based

attacks on keystroke based authentication [?], our approach does not require having knowledge of the console’s geometry. In other words, the size of the screen or the position of the pattern grid on the screen does not affect the accuracy of our attack. We also assume the video does not need to capture any content displayed on the screen. This assumption makes previous video-based attacks on pattern lock [?] inapplicable.

III. OVERVIEW OF OUR ATTACK

In this section, we give an overview of our attacking system which analyzes the user’s fingertip movement to infer the locking pattern. The system takes in a video segment that records the entire unlocking process. It produces a small number of candidate patterns to be tested on the target device. Figure 2 depicts the five steps of our attack:

① Filming and Video Preprocessing: The attack begins from filming how the pattern is drawn. The video footage can be filmed at a distance of about 2 meters from the user using a mobile phone camera (or 9 meters using a digital single reflex camera). After recording, the attacker needs to cut out a video segment that contains the entire unlocking process. We have shown that it is possible to automatically identify this video segments in some scenarios (Section IV-A). After cutting out the video segment, the attacker is then asked to mark two areas of interest from a video frame: one area consists of the fingertip used to draw the pattern, and the other consists of part of the device (see Figure 2 (b)).

② Track Fingertip Locations: Once the areas of interest are highlighted, a computer vision algorithm will be applied to locate the fingertip from each video frame (Section IV-B2). The algorithm aggregates the successfully tracked fingertip locations to produce a fingertip movement trajectory. This is illustrated in Figure 2 (c). Keep in mind that at this stage the tracked trajectory is presented from the camera’s perspective.

③ Filming Angle Transformation: This step transforms the tracked fingertip locations from the camera’s perspective to the user’s. We use an edge detection algorithm to automatically calculate the filming angle which is then used to perform the transformation (Section IV-C). For example, Figure 2 (c) will be transformed to Figure 2 (d) to obtain a fingertip movement trajectory from the user’s perspective.

④ Identify and Rank Candidate Patterns: In this step, our software automatically maps the tracked fingertip movement trajectory to a number of candidate patterns (Section IV-D).

We rank the candidate patterns based on a heuristic described in Section IV-D2. For instance, the fingertip movement trajectory in Figure 2 (d) could be mapped to a number of candidate patterns shown in Figure 11. Our approach is able to reject most patterns to leave no more than five candidate patterns to be tried out on the target device.

(5) Test Candidate Patterns: In this final step, the attacker tests the candidate patterns on the target device.

IV. IMPLEMENTATION DETAILS

A. Video preprocessing

The first step of our attack is to identify the unlocking process from the video footage. While all our participants (see Section V-A) consider this as a straightforward manual task, we developed a *simple yet effective* heuristic to automatically detect the video segment in some typical scenarios. Our heuristic is based on the following observations: (1) before or after unlocking, users often pause for a few seconds; (2) two consecutive on-screen operations (e.g. swiping, zooming etc.) typically expose some spatial-temporal motion characteristics. To test our hypothesis, we have recorded 50 video streams (where each video stream lasts around 2 minutes) of how ten of our participants drew patterns. During video recording, our participants firstly performed some on-screen activities such as web browsing and gaming as they wished; they then randomly opened up a pattern lock screen to draw a pattern and continued to perform other on-screen operations afterwards. We then analyzed frames that are associated with pattern drawing and those are not.

Figure 3 shows that all our participants paused at least 1.5 seconds before or after pattern drawing due to delay of the user or the device. We also found that identical on-screen activities often follow closely. For example, on several occasions our participants had to swipe several times to locate a program from the application list. These consecutive on-screen operations have some spatial-temporal motion characteristics that are different from pattern drawing. Figure 4 shows the spatial-temporal motion structure for two gestures, swiping and zooming, when they are performed once (a, c, e) and twice (b, d, f). This diagram indicates that the spatial-temporal motion of two identical on-screen activities contains one or more loop structures for which pattern drawing does not have.

Our heuristic is described in Algorithm 1. The input to the algorithm is a video capturing the unlocking process, and the output of the algorithm is a time-stamp tuple, $\langle \text{start}, \text{end} \rangle$, which marks the start and the end of a video segment. To locate the video segment of pattern drawing, we first filter out on-screen activities where the fingertip location does not change within a timeframe of 1.5 seconds (lines 4 and 11). This allows us to exclude some basic on-screen activities, such as clicking. We use the number of video frames, frameCount , as a proxy to estimate the time interval. Here, a time interval of 1.5s translates to 45 frames or 90 frames when the video was shot at 30 or 60 frames per second (FPS) respectively. We also use the spatial-temporal characteristics described above to exclude two consecutive swiping or zooming gestures (line 8). Finally, we exploit the observation that users typically paused

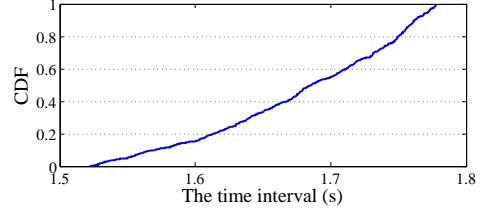


Figure 3. The cumulative distribution function (CDF) of the time interval between pattern drawing and other on-screen activities.

at least 1.5s before or after unlocking to locate the start and end points of pattern drawing (line 19).

Limitations Our heuristic is likely to fail if the user was typing using a Swype-like method (i.e. entering words by sliding a finger from the first letter of a word to its last letter) during video recording. In this case, our method will identify multiple video segments of which one may contain the pattern unlock process. If multiple segments are detected, the algorithm will ask the user to confirm which video segment to use. In this scenario, the first identified segment is likely to be the correct one. In practice, an experienced attacker would wait patiently to avoid this complicated situation by finding the right time for filming (e.g. for a screen lock, the time is just after the device is retrieved). The attacker could also watch the video to manually cut it to ensure the obtain the correct video segment. It is worthwhile to mention that automatically identifying the pattern unlocking process is not central to our attack because an experienced attacker can watch the video to manually cut it to ensure the tracking algorithm (described in the section) receives a quality input. Despite its limitations, our algorithm can reduce the efforts involved in some common scenarios.

B. Track fingertip locations

After cutting out the video segment of pattern drawing, we need to track the finger motions from the video segment. We achieve this by employing a video tracking algorithm called *Tracking-Learning-Detection (TLD)* [?]. This algorithm automatically detects objects defined by a boundary box. In our case, the objects to be tracked are the user's fingertip and an area of the device. These are supplied to the algorithm by simply highlighting two areas on the first frame of the video segment (see Figure 2 b). The algorithm tries to localize the fingertip from each video frame and aggregates the successfully tracked locations to produce a fingertip movement trajectory as an output (see Figure 2 c).

1) Generate The Fingertip Movement Trajectory: The TLD algorithm automatically detects objects based on the examples seen from the first frame. For each tracked object, the algorithm generates a confidence between 0 and 1. A tracking is considered to be successfully if the confidence is greater than a threshold. In our experiments, this threshold is set to 0.5 since this value is found to give good performance in our initial design experiments using 20 patterns². TLD has three modules: (1) a tracker that follows objects across consecutive

²To provide a fair evaluation, the patterns used in our initial test runs in the design phase are different from the ones used later in evaluation.

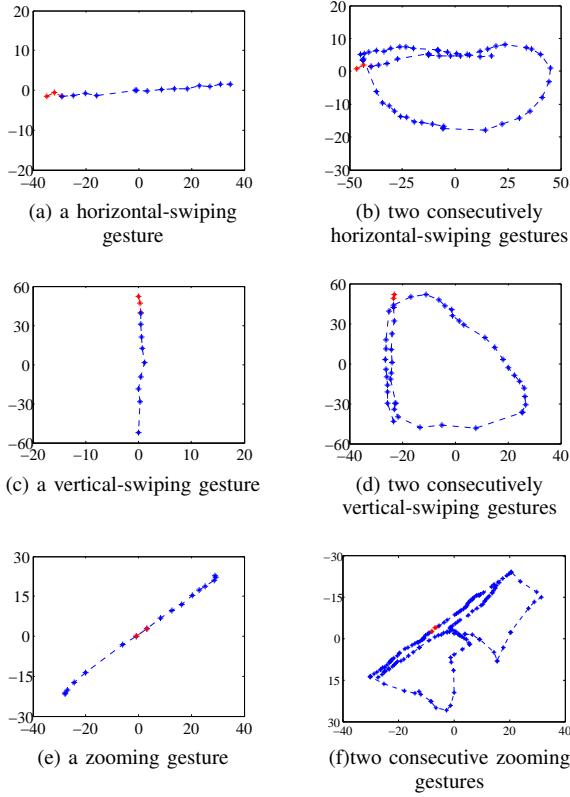


Figure 4. Spatial-temporal characteristics for performing an on-screen gesture once (a, c, e) and twice (b, d, f).

Algorithm 1 Unlocking process identification heuristic

Input:

IV: Video footage

frameCount: Pause threshold before or after unlocking

Output:

```

<start,end>: Start and end of the unlocking video segment
1: frames[] ← getVideoFrames(IV)
2: LEN ← getFramesLen(frames[])
3: for i = 1 : LEN - frameCount do
4:   sL ← hasFingertipChanged(frames[i : i + frameCount])
5:   if !sL then
6:     sNo = i + frameCount
7:     for j = sNo : LEN do
8:       if checkLoop(frames[j : LEN]) then
9:         eNo = i
10:        break;
11:       else if !hasFingertipChanged(frames[j : j + frameCount]) then
12:         eNo = i
13:         break;
14:       end if
15:     end for
16:     break;
17:   end if
18: end for
19: < start, end >← getTargetVideo(frames[], sNo, eNo)

```

frames under the assumption that the frame-to-frame motion is limited and objects are visible; (2) a detector to fully scan each individual frame to localize all appearances of the objects that have been seen in previous frames; and (3) a learner that estimates errors of the detector and updates the detector to avoid these errors.

The TLD learner automatically extracts features from the area of interest to build a K-Nearest Neighbor classifier [?] which is a part of the detector. In the following frames, the learner estimates the detection errors and generates new training examples (i.e. new appearances of the object) arose from object motion to re-train the classifier to avoid these errors. For each video frame, TLD calculates the tracking confidence and if the confidence is lower than the predefined threshold, the result of this particular frame will be discarded. This allows the algorithm to tolerate a certain degree of detection errors. Finally, the successfully detected object locations will be put onto a single image as the output. Detailed discussion of TLD can be found at [?]. Sometimes the algorithm may fail to detect the objects in many video frames due to poor selections of interesting areas. If this happens, our system will ask the user to re-select the areas to track. We have also extended TLD to report when a fingertip position is seen on the footage. This temporal information is recorded as the number of video frames seen so far. This is used to separate two possibly overlapping line segments described in Section IV-D.

2) *Camera Shake Calibration*: By default, the TLD algorithm reports the position of a tracked object with respect to the top-left pixel of the video frame. However, videos recorded by a hand-held device is not always perfectly steady due to camera shake. As a result, the top-left pixel of a video frame may appear in a different location in later frames. This can drastically affect the precision of fingertip localization, leading to misidentification of patterns.

Our approach to cancel camera shake is to record the fingertip location with respect to a fixed point of the target device. To do so, we track two areas from each video frame. One area is an edge of the device and the other is the fingertip. Both areas are highlighted on the first frame by the user. The location of a successfully tracked fingertip is reported as the relative coordinates of the two center points of the marked areas. This approach can also be used to calibrate the slight movements of the target device when drawing the pattern.

Example: To illustrate how our camera-shake calibration method works, considering Figure 5 where two areas are firstly marked by two bounding boxes in subfigure (a). Both areas will then be automatically detected by the TLD algorithm in following video frames as shown in subfigures (b) and (c). The coordinates of the two center points of each box are the values of x and y , and their relative positions are represented by ΔX and ΔY . For each frame where both areas are successfully tracked, we compute the relative coordinates, $(\Delta X, \Delta Y)$, which are reported as the location of the tracked fingertip.

Figure 6 shows the results when using TLD to process a video that was filmed with some camera shake effects. Here we show the results without (Figure 6 a) and with (Figure 6 b) camera-shake calibration. To aid clarity, we have converted the trajectories into the user's perspective. Without camera-shake

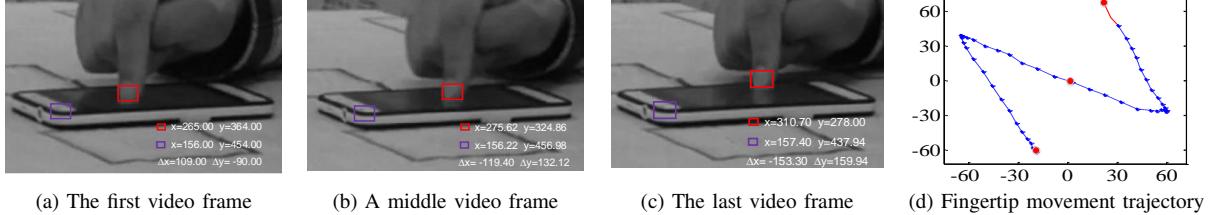


Figure 5. Tracking the fingertip movement trajectory. For each video frame, the system tracks two areas: one surrounds the fingertip and the other covers the edge of the device. The fingertip position is determined by computing the relative coordinates of the central points of the two areas. The red points highlighted in the final results (d) are the touching points tracked from the three video frames.

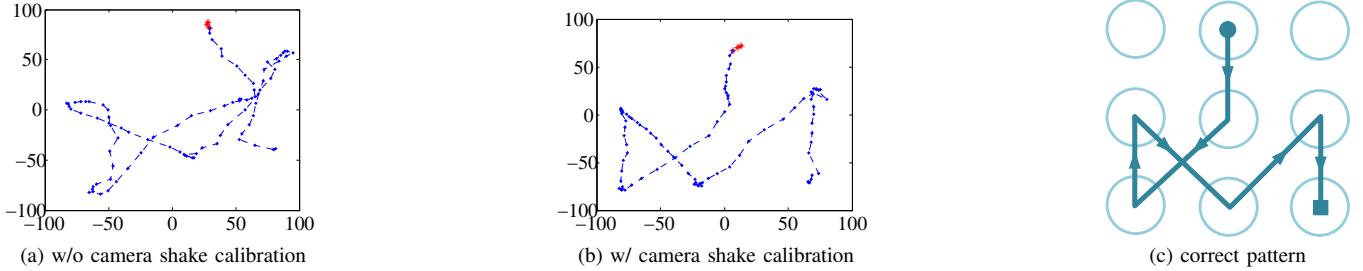


Figure 6. The resulted fingertip movement trajectories without (a) and with (b) camera-shake calibration. The correct pattern is shown in (c). To aid clarity we have transformed (a) and (b) to the user's perspective.

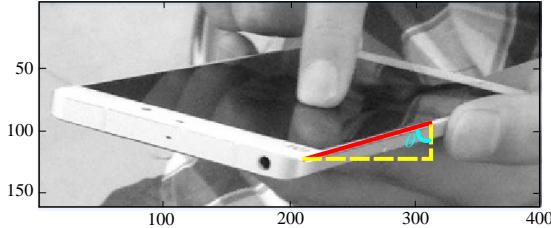


Figure 7. Filming angle calculation. The filming angle, θ , is the angle between the edge line of the device and a vertical line.

calibration, the resulted trajectory is significantly different from the actual pattern shown in Figure 6 (c). Because of this great difference, using Figure 6 (a) will lead to misidentification of candidate patterns. By contrast, processing the same video with camera-shake calibration produces Figure 6 (b). Clearly, this figure is more alike the correct pattern.

3) Noisy Points Calibration: During tracking process, the TLD algorithm may report mistaken position of a tracked object as rapid deformation of the tracked object. This can affect the shape of fingertip movement, leading to extract uncorrect geometric information of fingertip movement.

C. Filming angle transformation

In practice, the filming camera will not directly face the target device to avoid the filming activity to be noticed by the user. As a result, the fingertip movement trajectory generated by the tracking algorithm will look differently from the actual pattern. For example, for the pattern presented in Figure 2 (a), if the video is filmed from the attacker's front-left to the target device (i.e. with a filming angle of approximate 45 degrees), we get the trajectory shown in Figure 2 (c). Using this trajectory without any postprocessing will lead to

misidentification of candidate patterns. Therefore, we must transform the resulted trajectory to the user's view point. To do so, we need to estimate the angle between the filming camera and the target device. Our approach is described as follows.

We use an edge detection algorithm called Line Segment Detector (LSD) [?] to detect the longer edge of the device. The filming angle is the angle between the detected edge line and a vertical line. This is illustrated in Figure 7. In Section VI-E, we show that a slight estimation error of the filming angle has little impact on the success of the attack. By default, we assume that the pattern grid is presented in the portrait mode³. If this is not the case, i.e. the pattern grid is shown in the landscape mode, we need to use the shorter edge of the device to calculate the filming angle. We believe that an attacker interested in a particular target would have some knowledge of how the pattern grid of the target device is presented under different orientation modes and be able to identify the device orientation by watching the video. There are also other methods to be used to identify the filming angle [?].

Based on the estimated filming angle, θ , we use the following formula to transform the tracked fingertip movement trajectory from the camera's view point to the user's:

$$S = TS' \quad , \quad T = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad (1)$$

where T is a Transformation Matrix, S' is the coordinate of a point of the tracked trajectory, and S is the resulted coordinate after the transformation. For each video frame, our algorithm individually calculates the filming angle and perform the transformation, because the filming angle may change across video frames.

³The pattern grid of the Android native pattern lock is always presented in the portrait mode regardless of the orientation of the device.

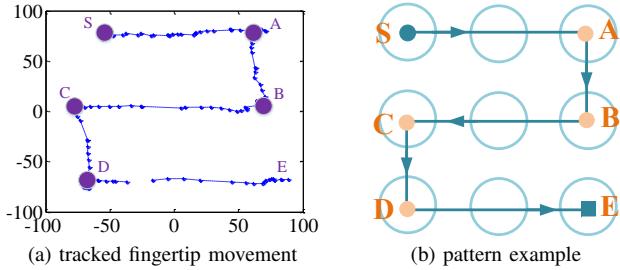


Figure 8. This figure shows the tracked fingertip movement trajectory (a) of a pattern (b). Point S on (a) is the starting point and points A, B, C, and D on (b) represent four turning points.

Algorithm 2 Line Segment Identification

Input:

struct $T[]$: Temporal information of each tracked location
 $timeTh$: Threshold of whether two line segments are overlapping

Output:

$tp[]$ Turning points of fingertip movement.

```

1: for each fingertip movement with temporal sequences  $T[]$ 
  do
2:    $tpNum = 0;$ 
3:   struct  $lines[] \leftarrow getLines(T[])$ 
4:    $lNum \leftarrow getLinesNumber(lines[])$ 
5:   for  $i = 1 : lNum$  do
6:     if  $checkOverlap(lines[i], timeTh)$  then
7:        $p[tpNum + +] \leftarrow getOverlapPoints(line[i])$ 
8:     end if
9:      $p[tpNum + +] \leftarrow getTurningPoints(line[i])$ 
10:   end for
11: end for
12:  $tp[] = p[0 : end - 1]$ 
```

D. Identify and rank candidate patterns

In this step, the fingertip movement trajectory will be mapped to a number of candidate patterns to be tested on the target device. Our goal in this step is to exclude as many patterns as possible and only leave the most-likely patterns to be tried out on the target device. Our approach is to use the geometry information of the fingertip movement trajectory, i.e. the length and direction of line segments and the number of turning points, to reject patterns that do not satisfy certain criteria. In this section, we first describe how to identify overlapping line segments and extract length and direction information before presenting how to use the extracted information to identify and rank candidate patterns.

1) *Extracting Structure Information:* A pattern can be defined as a collection of line segments where each line segment has two properties: the length of the line, l , and the direction of the line, d . Figure 10 shows all the possible directions of a 3×3 pattern grid. We define a pattern, P , as a collection of line segment properties, $P = \{L, D\}$. Here $L = \{l_1, l_2, \dots, l_n\}$ is a collection of the lengths of all line segments (that are numbered from 1 to n) of the pattern, and $D = \{d_1, d_2, \dots, d_n\}$ is the collection of directions for all line

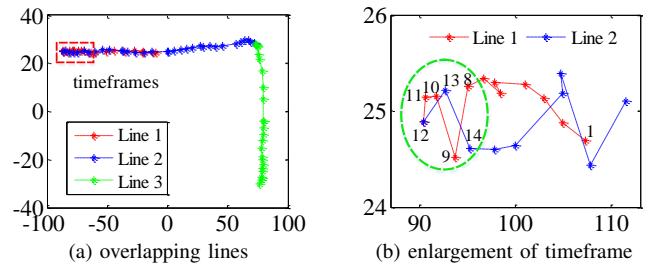


Figure 9. Separating two overlapping line segments by checking the number of overlapping points within a timeframe.

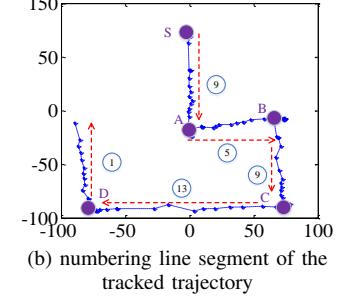
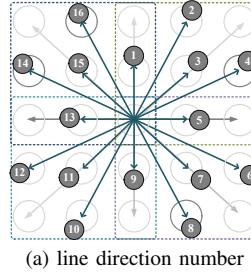


Figure 10. All possible line directions for a 3×3 Android pattern grid.

segments in L . Algorithm 3 describes how P is extracted. We extract the length and the direction of each line segment from the tracked fingertip movement trajectory and store them into arrays $L[]$ and $D[]$ respectively.

Identify Line Segments The first step of geometry information extraction is to identify individual line segments from the trajectory. This can be achieved by finding turning points, the start and the end points of the pattern, because two points define a line segment. For example, turning points, A and B, in Figure 8 defines a line segment, AB. In Algorithm 2, we use a linear fitting method [?] to discover turning points (line 3). A specific challenge here is how to separate two overlapping line segments (see Figure 12 c for an example). It is to note that up to two lines can be overlapped on a pattern grid. The naive linear fitting algorithm would consider two overlapping segments to be a single line as their points stay close to each other. We overcome this problem by using the temporal information (that is recorded by the tracking algorithm) to separate two overlapping points. To do so, we visit all tracked points of each line segment given by the linear fitting algorithm (line 5) within a timeframe ($timeTh$) of 20 video frames for a video of 30 FPS (40 for a video of 60 FPS). For each point, we calculate its Euclidean distances to all other points within the timeframe. We consider two points to be overlapping if their distance is less than 5 pixels. For a video shot at 30 FPS, we consider there exist two overlapping line segments if 5 (10 for a 60 FPS video) or more overlapping points in the timeframe. Again, these threshold values were determined through our initial design experiments. Finally, we consider the center of all points as the turning point of the two overlapping line segments and use turning point to separate the two lines. *Example:* As an example, consider a fingertip movement

Table I
MAPPINGS FROM LINE SLOPES AND FINGERTIP-HORIZONTAL MOVEMENTS
TO DIRECTION NUMBERS

Direction No.	1	2	3	4	5	6	7	8
slope ($L \rightarrow R$)	$+\infty$	2	1	$\frac{1}{2}$	0	$-\frac{1}{2}$	-1	-2
Direction No.	9	10	11	12	13	14	15	16
slope ($R \rightarrow L$)	$-\infty$	2	1	$\frac{1}{2}$	0	$-\frac{1}{2}$	-1	-2

Algorithm 3 Candidate Pattern Identification Algorithm

Input:

$L[]$: Relative line length

$D[]$: Direction number (see Figure 10)

tn : Number of turning points of fingertip trajectory

$lengthTh$: Threshold of considering two lines to have the same length

$directionTh$: Threshold of considering two lines to be in the same direction

Output:

$P[]$: Candidate patterns

```

1: for each possible pattern  $p$  with  $tn$  turning points do
2:    $n \leftarrow getLineNumber(P[])$ 
3:    $pL[] \leftarrow getRelativeLength(p)$ 
   /*Relative line length for pattern  $p$ */
4:    $pD[] \leftarrow getDirection(p)$ 
5:   if  $match(pL[], L[], lengthTh)$  then
6:     if  $match(pD[], D[], directionTh)$  then
7:        $P[] \leftarrow p$ 
8:     end if
9:   end if
10: end for
11:  $P[] \leftarrow sort(P[])$ 

```

trajectory shown in Figure 9 (a). The red rectangle on the figure is a timeframe consisting of 20 tracked points. If we zoom in on the timeframe, we get Figure 9 (b) where a point is labelled with a frame number according to when the point was seen, starting from 1 for the earliest point. In this example, there are more than 6 overlapping points in the timeframe, which are marked by a green circle. We use the center point (No.10) of the overlapping points as the turning point to separate the two line segments.

Extract the Line Length The physical length of a line segment depends on a number of device-specific parameters: the sizes of the screen and the pattern grid, and the space between two touch dots. To ensure our approach is independent of the device, we normalize the physical length of a line segment to the shortest line found on the tracked trajectory. For the example shown in Figure 8 (a), the line lengths for segments, SA, AB, BC, CD, and DE, are $2l_s, l_s, 2l_s, l, 2l_s$, respectively. Here segments AB and CD have the shortest length, l_s . The physical length of a line segment is calculated by computing the Euclidean distance between the start and the end points of a segment.

Extract Direction Information In addition to the line length, we also want to know to which direction the finger moves. This information is useful for inferring which dots are selected

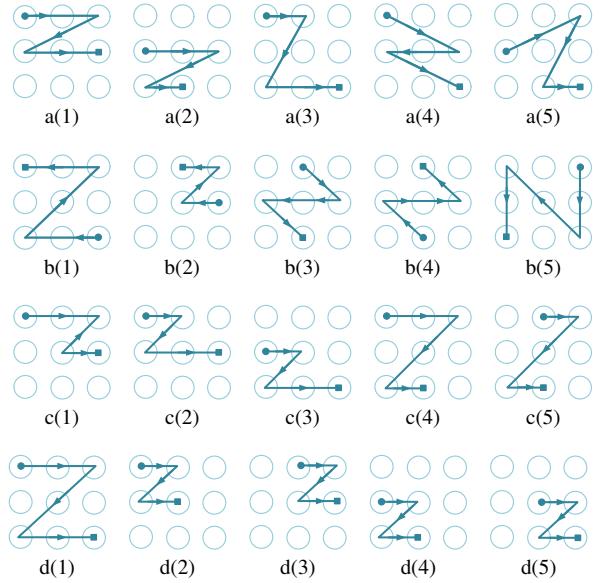


Figure 11. Possible mappings for the tracked fingertip movement trajectory presented in Figure 2 (d).

to unlock the pattern. Figure 10 (a) shows all possible 16 directions on a 3×3 pattern grid. The directions are numbered from 1 to 16 in clockwise. For each line segment of the tracked trajectory, we calculate its line slope and the horizontal movement of the finger (i.e. left \rightarrow right or vice versa). This information will then be checked against Table I to determine the direction number of the line segment. The horizontal movement of the fingertip is determined by first using the temporal information to find out the start and the end points of the line and then comparing the horizontal coordinates of the two points. The line slope is also computed based on the coordinates of the start and the end points of the line segment. Figure 10 (b) gives the direction number of each tracked line segment of a fingertip movement trajectory.

2) *Map the Tracked Trajectory to Candidate Patterns:* In this step, we use the extracted geometry information to map the fingertip movement trajectory to a small number of candidate patterns which will then be ranked using a heuristic. This process is described in Algorithm 3.

Identify Candidate Patterns Our implementation simply enumerates all possible patterns for a given pattern grid to identify candidate patterns, starting from the top-left touch point. We reject patterns that do not meet the requirements that the correct pattern is expected to have. The requirements are the number of line segments (this is checked by counting the number of turning points), and the length and the direction for each line segment. This is an *automatic* process performed by our software system without any user involvement. We consider two line segments having the same length and slope if the difference between them is less than a threshold. Specifically, the relative length threshold, $lengthTh$, is set to 1.12 and the slope threshold, $directionTh$, is set to 0.25. To determine the thresholds, we have evaluated a range of possible values in our initial design experiments to chose the best performing values.

Example: We use the pattern depicted in Figure 2 as an

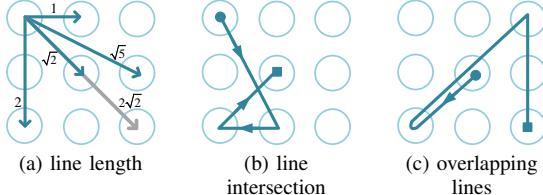


Figure 12. Illustrations of the terminologies used in Equation 2.

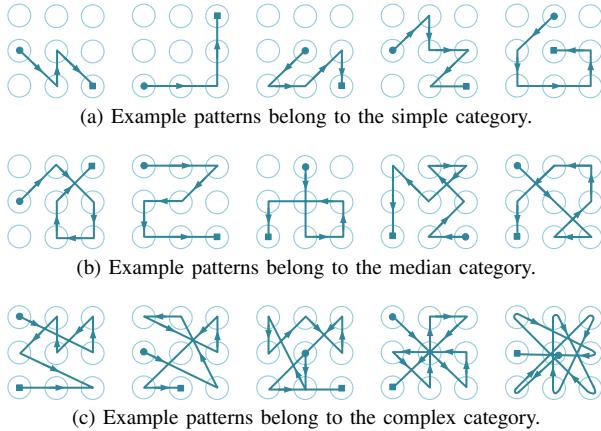


Figure 13. Examples of patterns collected from our participants. Patterns are grouped into *simple*, *median* and *complex* categories, according to their complexity scores.

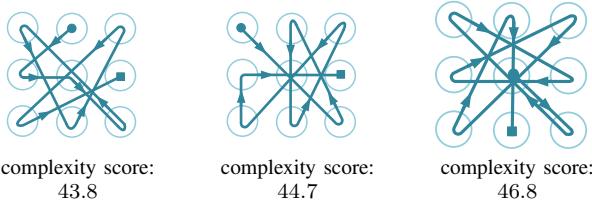


Figure 14. Three most complex patterns on a 3×3 grid based on Equation 2.

example to describe our algorithm. Figure 11 gives several possible mappings for the fingertip movement trajectory shown in Figure 2 (d). For this particular trajectory, the collections of lengths and directions are $L = \{l, \sqrt{2}l, l\}$ and $D = \{5, 11, 5\}$ respectively. Any pattern that does not meet L or D should not be considered as a candidate pattern for this trajectory. For this reason, Figure 11 a(1)–a(5) will be rejected. Take Figure 11 a(1) as an example, the line lengths and directions for all four line segments of this pattern are $\{l, \frac{\sqrt{5}}{2}l, l\}$ and $\{5, 12, 5\}$ respectively. It does not meet the expected L or D and should be rejected. The patterns presented in b(1)–b(5) and c(1)–c(5) of Figure 11 will also be rejected for the same reason.

Rank Patterns Candidates patterns are then ranked using a simple heuristic. The heuristic assumes a pattern starting from left dot of the grid is more likely to be the correct pattern over a pattern starting from a right dot. This assumption is supported by recent studies which show that people tend to select a left dot as the starting point to construct a pattern [?]. If two candidate patterns start from the same dot, we consider the pattern with a longer total line length is more likely to be the correct pattern. Using these criteria, the five

candidate patterns are ranked in order from subfigures d(1) to d(5) in Figure 11. Therefore, an attacker would first try the candidate pattern presented in Figure 11 d(1). This attempt will lead to a successful attack for the example presented in Figure 2. Our experimental results show that this heuristic is effective.

V. EXPERIMENTAL SETUP

A. Data Collection

The patterns used in our evaluation were collected from users who use at least one Android device (a smartphone or a tablet) on a daily basis. To collect the patterns, we have distributed over 1,000 survey forms and collected back 215 valid forms, resulting in 120 unique patterns⁴. Our participants include 95 females and 120 males who were undergraduate or postgraduate students in our institution. The majority of our participants are in an age group of under 30.

To collect the patterns, we have conducted a “pen-and-paper” survey by asking participants to fill in an anonymized questionnaire. The questionnaire and survey were approved by the research ethics board (REB) of our institution. We have made sure that our survey complied with strict privacy regulations. For example, we did not collect any personally identifiable information other than the gender and age group of the participant. Our participants were well informed on the purpose of the study and how the data will be managed and used. The survey forms were distributed as voluntary homework so that the participants can take the survey form away to fill in. Users were invited to return the survey form anonymously within three weeks to a dedicated, locked mailbox, if they wish to participate in the study. To avoid a user submits multiple copies of the same form, each survey form is given a unique, randomly generated 32-digit number.

Overall, 37.6% of our participants confirmed that they use pattern lock as the screen lock to protect their Android devices on a daily basis; and 33% of those who do not use a pattern as their screen lock said that they are often required to use a pattern for authentication by an application like Alipay. Furthermore, 60% of our participants also indicated that the pattern they provided is currently being used or have been used in the past by themselves. Other participants (often those did not use a locking pattern on a daily basis) indicated that they have provided a pattern which they would like to use if a locking pattern is required. Based on this information, we are confident that our patterns represent some of the real world patterns. Finally, all participants believe that a complex pattern provides stronger protection than a simple one.

B. Pattern Complexity Classification

We quantify the complexity of a pattern using the complexity (strength) score proposed in [?]. The complexity score, CSP , of a pattern, P , is defined as:

$$CSP = S_P \times \log_2(L_P + I_P + O_P) \quad (2)$$

⁴Patterns are available at: <https://github.com/yeguixin/patterns>

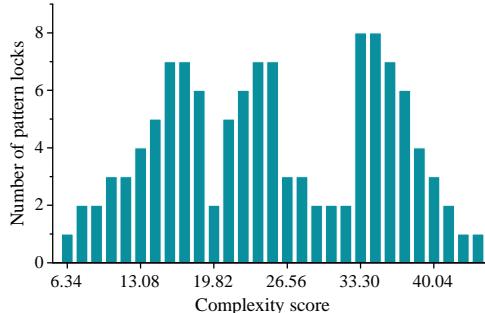


Figure 15. The distribution of complexity scores for the patterns given by our participants.

Table II
SCREEN SIZES FOR THE TEST PHONES

Size	Brands	MI4	Honor7	Note4
Height(cm) × Width(cm)		13.9 × 6.9	14.3 × 7.2	15.4 × 7.9

where S_P is the number of connected dots, L_P is the total length of all line segments that form the pattern (see Figure 12 a), I_P are the number of intersections (which are also termed as “knight moves” in some prior work [?], see Figure 12 b) and O_P are the number of overlapping linear segments (see Figure 12 c). To calculate the line length, we assume the length between two horizontally or vertically adjunct dots is one. Thus, our method is independent of the size of the screen and the grid.

Intuitively, the more connected dots (S_P), line segments (L_P), intersections (I_P) and overlapping line segments (O_P) that a pattern has, the more complex it is. For example, the patterns shown in Figure 13 (c) use all the nine dots of the grid, and have at least seven line segments and three intersections.

Based on the complexity score, we divide the collected patterns into three complexity categories: *simple*, *median* and *complex*. A simple pattern has a score of less than 19, a median complex pattern has a score between 19 and 33, and a complex pattern must have a score greater than 33. This classification gives us roughly 40 patterns per category. Figure 13 gives some examples for each category while Figure 15 shows the distribution of these patterns according to their complexity scores. Based on this definition, the most complex pattern on a 3×3 grid has a score of 46.8 (see Figure 14). The complex scores of the patterns we collected range from 6.4 to 46.8.

C. Video Recording and Preprocessing

Recording Devices We used three smartphones for video recording: an Apple iPhone4S, a Xiaomi MI4 and a Meizu2 phones. Each mobile phone was used to record 40 patterns with a 1080p HD resolution of 30 FPS under different settings.

User Participation We recruited ten postgraduate students (five male and five female students) from our institution to reproduce the 120 patterns (collected from the users) and the 60 most complex patterns (see Section VI-A) on three target mobile phones: a Xiaomi MI4, a Huawei Honor7 and a Samsung Note4 smartphones. Table II lists the screen size for each target mobile phone.

Video Recording Setup By default, we used the Android 3×3 native pattern grid, but we evaluated our approach using other pattern grids with different sizes in Section VI-H. We recorded each pattern under three filming angles, 45, 90 and 135 degrees, by placing the camera on the left-front, front, and right-front of the target device respectively. By default, the video was recorded indoor during daytime under a natural lighting condition. In Section VI-D we evaluated our approach under different lighting conditions both indoor and outdoor. By default, videos were recorded at a distance of 2 meters from the target device and we evaluated the impact of the filming distance in Section VI-H.

Video Filming Before recording, our participants were given the opportunity to practice a pattern several times, so that they can draw the pattern at their natural speed. On average, this practice session took 10 trials per user per pattern. When drawing the pattern, some participants sat, while others stood, some hold the device by hands, while others placed it on a table. Each pattern was drawn on three target devices and recorded under three filming angles. Thus, for the 120 patterns collected from users, we recorded 1,080 videos in total.

Video Preprocessing For each video stream, we used the algorithm described in Section IV-A to cut out the video segment of the unlocking process. We left around 200 to 300 milliseconds of the video segment before and after the pattern unlocking process. To track the fingertip locations, we used Windows Movie Maker to highlight two areas of interest on the first frame of the video segment: one area surrounds the fingertip, and the other contains an edge of the phone (see Section IV-B2).

Implementation Our prototyped attacking system built upon a TLD library [?]. The developed software ran on an Intel Core i5 PC with 8GB RAM. The operating system is Windows 10. Our implementation can be ported onto Android or Apple iOS systems, which is our future work. On our evaluation platform, our software takes less than 30 seconds to process a video to produce candidate patterns.

VI. EXPERIMENTAL RESULTS

In this section, we first present the overall success rate for cracking the 120 patterns collected from our participants plus the top 60 most complex patterns on a 3×3 pattern grid. Our results show that our approach can successfully crack over 95% of the patterns using no more than five attempts. We then analyze how the success rate is affected by the filming distance, filming angles and camera shake. Finally, we demonstrate that direct observations lead to poor performance before evaluating our approach on alternative pattern grids.

A. Overall Success Rate

Result 1: *We can successfully crack over 95% of the patterns in five attempts and complex patterns are less secure compared to simple patterns under our attack.*

In this experiment, videos were recorded from a distance of 2 meters away from the target device. This mimics a scenario where the adversary sits at the next table to the user in a public space (e.g. a restaurant). The smartphones used for

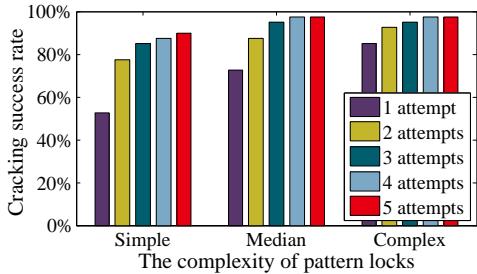


Figure 16. For each pattern category, the figure shows the success rate using no more than 1, 2, 3, 4 and 5 attempts.

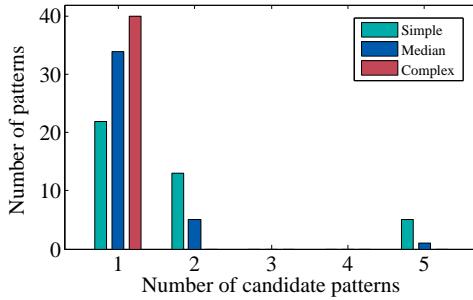


Figure 17. The distribution of candidate patterns for each category. No more than 5 candidate patterns were generated by our algorithm.

Table III
TRACKING PRECISION VS FILMING DISTANCE

Distance	1 m	2 m	3 m	3.5 m
fingertip	100%	98.7%	80.9%	68%
device edge	100%	99.4%	90.6%	69%

filming in this experiment were hand-held. Figure 16 shows the success rate for cracking different types of patterns within 1, 2, 3, 4 and 5 attempts. For all the patterns used in this evaluation, our approach does not generate more than five candidate patterns. For complex patterns, we are able to crack all except one in the first attempt. For simple and median patterns, the success rate increases with more tries. In one attempt, we are able to successfully crack 60% and 87.5% of the simple and median patterns respectively. With two attempts, the success rate increases to 87.5%, and 95% for simple and median patterns respectively. Using five attempts, we are able to crack all simple patterns and all but one median patterns. The reason that we failed on one median and one complex patterns is because of some blur motions of the video footage (probably caused by the video compressing algorithm), which leads to many tracking failures. But we are able to crack the same pattern using a video filmed by a different device. It is important to note that the Android OS will not lock the device unless seeing more than five failed tries [?]. This means, in practice, our approach is able to successfully crack most locking patterns.

Another interesting observation is that in contrast to many people's intuition, complex patterns do not provide stronger protection under our attack because most of these patterns can be cracked in one attempt. This is because although complex patterns can better protect the user against direct observation techniques like shoulder surfing [?], their unique graphical structures often allow our system to reject most patterns to

produce a single candidate pattern. This is confirmed by Figure 17. It shows that for most median and all complex patterns, our system produces one candidate pattern – the correct one for most of our test cases.

We also evaluated our approach using the top 60 most complex patterns (according to Equation 2) on a 3×3 grid. To evaluate our approach on a wide range of patterns, we exclude patterns that are simply a rotation to an already chosen pattern. Figure 14 illustrates three highly complex patterns which have a complexity score between 43.8 and 46.8. The three patterns use all the nine dots of the grid and have a larger number of line segments, intersections and overlapping lines when compared to simpler patterns. Because of their complex graphical structures, it would be hard to remember these patterns using direct observation techniques. In this experiment, we can crack all the complex patterns in one attempt. This result reinforces our claim that complex patterns are less security under video-based attacks.

B. Impact of Filming Distances

Result 2: We can crack over 80% of the patterns in five attempts, if the video was filmed using a smartphone within a distance of 2.5 meters away from the target.

We would like to know how the filming distance affects the success rate of the attack. To do so, we have asked our participants to randomly select all 120 pattern locks and we varied the filming distance from 1 meter to 3.5 meters. Figure 19 shows how the cracking success rate changes as the filming distance increases. There is slight differences in the success rate between this diagram and Figure 16 because we used less patterns in this experiment. When the filming distance is less than 2 meters, our approach can crack all patterns in five attempts. The success rate drops significantly when the filming distance is greater than 2.5 meters. The quality of the video filmed by a mobile phone tends to drop significantly with many object deformations. The degradation of the video quality makes it difficult for the TLD algorithm to successfully track objects across video frames. This is confirmed by Table III which shows that the tracking precision for the fingertip and the device edge drops from around 99% to 68% when the filming distance increases from 2 meters to 3.5 meters. The increased tracking failures result in an increased number of missing points on the tracked trajectory, leading to a deteriorative performance in identifying candidate patterns. This can be seen from Figure 18 where the quality of tracking clearly decreases when the filming distance is greater than 3 meters. Nonetheless, our approach can achieve a high success rate when the filming distance within 2.5 meters. Such a distance allows an attacker to record the video without raising suspicions in many day-to-day scenarios (some of these are described in Figure 1).

We also evaluated our approach on videos filmed using a Nikon D90 single-lens reflex (SLR) camera with a 105mm lens. The SLR camera was placed from a distance of 9 meters away from the target device. For this set of videos, we are able to achieve the same performance when compared to using videos filmed by a mobile phone camera with a 2-meter

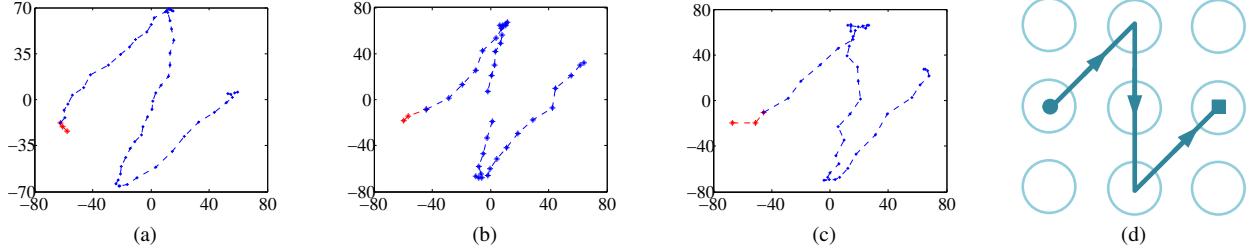


Figure 18. Tracked fingertip trajectories (user's perspective) for the pattern shown in (d) from a video filmed from a distance of 2m (a), 3m (b), and 3.5m (c) respectively away from the target device. The tracking quality decreases when the filming distance is greater than 3m.

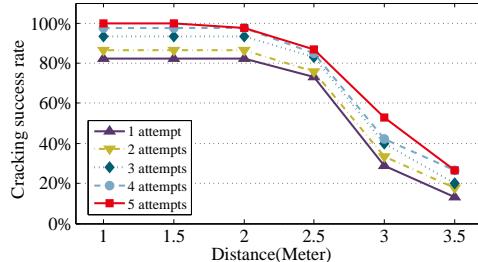


Figure 19. Impact of the filming distance.

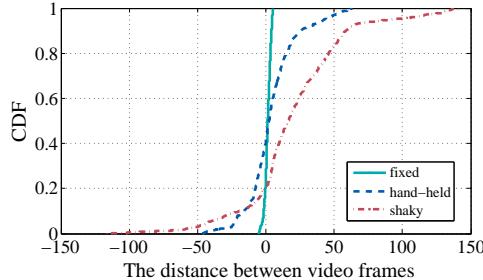


Figure 20. The cumulative distribution function (CDF) for different video recording modes.

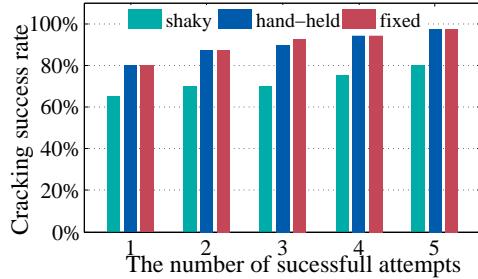


Figure 21. Impact of camera shake. Our approach has the same success rate under the hand-held and the fixed modes and the performance degradation under the shaky mode is modest.

filming distance. The further filming distance is largely due to better video quality brought by the advanced SLR camera and the lens. Therefore, in practice, an attacker can also use a professional video recording device to launch the attack from a further distance.

C. Impact of Camera Shake

Result 3: Our method can tolerate a certain degree of camera shake in the hand-held mode.

In this experiment, we used an iPhone4S smartphone to record how a pattern is drawn on a Huawei Honor7 phone. This experiment was carried out under three settings: *fixed*,

Table IV
LIGHTING CONDITIONS

Scenarios	Indoor	Indoor	Indoor	Outdoor
Time	nighttime	nighttime	daytime	daytime
Light Source	warm LED	white fluorescent	sunlight	sunlight
Light Intensity (Lux)	55 – 70	70 – 100	150–240	500–9500

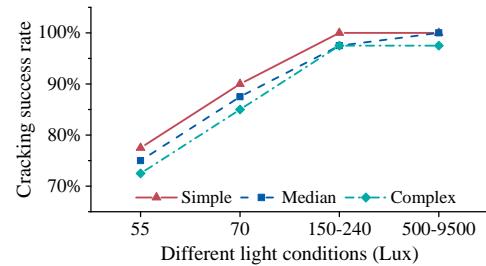


Figure 22. The cracking success rate within five attempts under different lighting conditions.

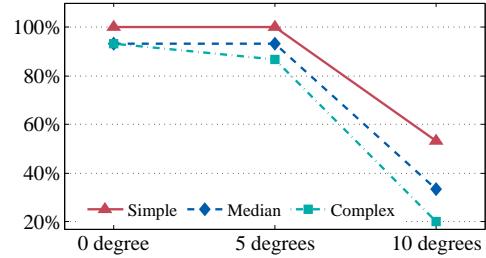


Figure 23. Impact of estimation errors of filming angles.

hand-held and *shaky*, where the filming device was respectively fixed using a tripod, hand-held, and hand-held but with constant movements of approximate 2cm in the horizontal or the vertical directions. The recording device was placed on the left-front, front, and right-front of the target device. In the experiment, we fixed the target device on a table using double-sided tapes.

We use a reference point to quantify camera shake. The point is the center position of an area of the target device. The area is marked by a boundary box on the first frame (see Figure 5). We calculate the difference (in terms of pixels) for where the reference point was seen in two consecutive video frames. We then use the difference to measure the degree of camera shake. Figure 20 shows the cumulative distribution function (CDF) of camera shake under the three different filming settings. Here, the wider the distribution is, the less steady the filming is. The 'shaky' mode is least stable where the difference of the reference point between two video frames

can be up to 250 pixels.

Figure 21 shows that our approach has the same performance under the hand-held and the fixed modes. The modest camera sake under the hand-held mode has little impact on performance thanks to our camera-shake-calibration method. We observe deteriorative performance under the shaky mode, but the performance degradation is modest (80% vs 97% in 5 attempts). In reality, an attacker would avoid drastic camera shake by firmly holding the video recording device.

D. Impact of Lighting Conditions

Result 4: *Low-light has a negative impact on the success rate of the attack but our approach can still break over 70% of the patterns when the video was filmed in a low-light environment.*

In this experiment, videos were recorded under different lighting conditions both indoor and outdoor. The experimental settings are given in Table IV. The light intensity of these condiditions ranges from 9500 lux (strong light), onto 240 lux (normal light), and 55-70 lux (low light). These represent some of the day-to-day scenarios where filming can take place. For each setting, we asked each of our 10 participants to draw all collected 120 patterns on a Xiaomi MI4 phone. We used an iPhone4S phone to record the video. The filming camera was place on the left-front, front, and the right-front of the target device from a distance of 2 meters.

Figure 22 shows that the success rate increases when video filming were performed in a brighter lighting condition as the light intensity changes from 55 lux to 9500 lux. This is expected as low-light leads to increased video noise, blurred motion and poor focus, which all have a negative impact on the TLD algorithm. Nonetheless, our attack can still crack over 70% of the patterns in a filming environment with low light.

E. Impact of Filming Angle Estimation

Result 5: *Our attack performs well when the error of filming angle estimation is less than 5 degrees.*

Our approach needs to transform the tracked fingertip movement trajectory to the user's perspective based on an estimation of the filming angle (Section IV-C). Because our filming angle estimation algorithm gives highly accurate estimation, we did not find the estimation error to be a problem in our experiments. Nonetheless, it is worth studying how the estimation error affects the success rate of our attack. To do so, we deliberately added an error of 5-10 degrees to the estimation.

Figure 23 shows the results of this experiment. When the error is less than ± 5 degrees, there is little impact on *complex* patterns and no impact at all on *simple* and *median* patterns. However, an estimation error of more than 10 degrees can significantly affect the success rate. Given such an error, the resulted trajectory after transformations will be significant different from the correct pattern. For example, when the estimation error is 10 degrees from the true value, on average, 0.8, 2.6 and 4.2 line segments per pattern respectively will be incorrectly labelled for *simple*, *median* and *complex* patterns. This explains why the success rate for complex patterns drops significantly when there is an error of 10 degrees.

Table V
TOUCH-SCREEN SIZES FOR THE TEST PHONES

Brands Size	iPhone4S	Honor7	Samsung Tab
Height(cm) \times Width(cm)	11.5 \times 5.9	14.3 \times 7.2	24.2 \times 15.0

Table VI
THE MAIN PARAMETERS FOR THE PHONE CAMERA

Brands Parm	iPhone6	Vivo X7	MI4	Note4
Frame Rate (fps)	30	30	30	30
Pixels	8mp	13mp	13mp	16mp
Focus (mm)	4.15	4	4	4.2
Sensitivity (ISO)	3200	3200	3000	5000

F. FIXED: Impact of Different Targets and Cameras

Result 6: *The size of smartphone touch-screen and brands of mobile camera had little effect on the success rate.*

Intuitively, the success rate may be influenced by the size of touch-screen because the length of pattern will be short as the screen size decreases. Likewise, the accuracy also may be influenced by the quality of video as the different quality of phone cameras. To evaluate this effect, we ask 10 participants to randomly select 60 patterns (30 patterns for each category). We use iPhone6, Xiaomi MI4, Vivo X7, Samsung Note4 to record unlock videos while the patterns are drawn on iPhone4S⁵, Huawei Honor7 and Samsung Galaxy Tab E. The distance is 2m away from the tester during record the unlocking process. Table V shows the sizes of screen vary from 3.5inch (iPhone4S) to 9.6inch (Samsung Tab) which involve almost all screen size of smart devices used in daily life. Table VI presents the mobile cameras and their main parameters used in this experiment. The parameters for the four mobile cameras have few differences and these differences determine the function of the cameras in specific shooting scenes. For example, Note4 camera have a longer focal length than other three mobile camera, which means that its effective shooting distance is longer than others. We believe that the four mobile cameras represent the most common cameras used by people of all walks of life.

Figure 24 (a) shows our approach has the same performance under different sizes of phone screens. We observe the success rate reach 98.3% within five attempts when the screen size is 3.5inch, the smallest screen on the mobile market, and our method can crack all tested pattern locks when the size of screen is larger than 5.2inch (Honor7). We also evaluated the our approach using different types of phone cameras for filming videos. The result is presented in Figure 24 (b). Our method can reconstruct all patterns within five trials using iPhone6 and Note4 cameras because they are able to capture high quality videos. Anyway, the success rate can reach 96.7% and 98.3% within five attempts respectively using MI4 and Vivo X7 cameras. This prove that our method perfects well with common used phone cameras.

⁵Participants draw patterns on the login interface of Alipay installed on the phone as Apple operating system do not support pattern lock authentication mechanism.

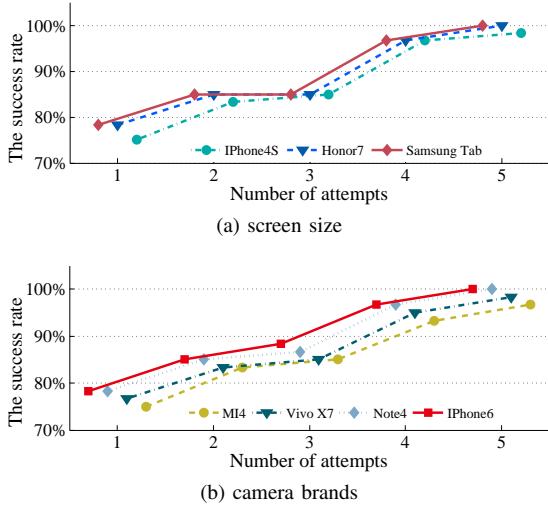


Figure 24. The cracking success rate under different screen size and brands of camera. The difference in horizontal axis is for the purpose of illustration.

G. Inferring Patterns with Eyes

Result 7: Our attacking methodology significantly outperform direct observation techniques.

In this experiment, we investigate whether an attacker can infer the pattern by simply watching the video or through direct observations. To answer this question, we asked each of our 10 participants to watch 60 videos (where a pattern was drawn by other participants) to guess the pattern. We only played the video segment during which a pattern is drawn to the participant (around 3 seconds per video). To familiarize participants with the process, we played five sample videos and showed the correct patterns at the end of each video to our participants before the experiment. Each participant then had 10 minutes to watch a video and five chances to guess a pattern. They could adjust the playing speed and replay the video multiple times as they wished.

Figure 25 (a) shows the success rate of pattern guessing with bare eyes. Our participants correctly guessed for nearly half of the simple patterns in five attempts. However, they found that it is difficult to infer complex patterns with many line segments, overlapping lines and intersections. The success rate of guessing complex patterns is less than 10% in five attempts. This is not a surprising result because although it is possible to correctly guess patterns with simple structures by watching the video, doing so for patterns with more complex structures is much harder.

We also asked participants to directly observe how a pattern was drawn from a distance of 2 meters away from the target device. The intuition behind this evaluation is that human eyes can catch richer information over a video camera. The results of this experiment are shown in Figure 25 (b). As can be seen from the diagram, although the success rate is improved compared to directly watching the video, the chances for guessing the correct pattern in 5 attempts are quite low. In fact, the success rates are just 48.3%, 38.3% and 11.7% respectively for simple, median and complex patterns.

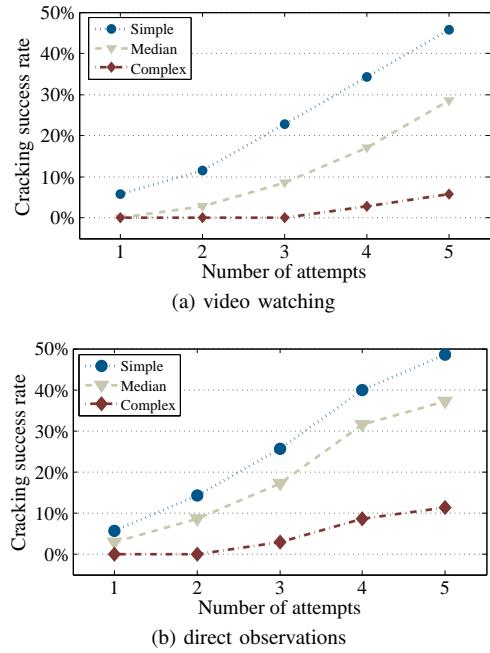


Figure 25. Success rates of guessing patterns through watching the video (a) or direct observations (b).

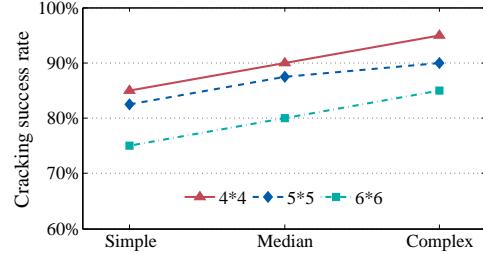


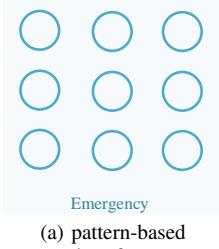
Figure 26. Success rates of our attack for different locking grids.

H. Evaluation on Other Pattern Grids

Result 8: A pattern grid with more dots provides stronger protection but our attack can still crack most of the patterns.

There are a few applications (such as CyanLock) and customized ROMs available to increase the size of the pattern grid from 3×3 to 4×4 , 5×5 , and 6×6 . Although a 3×3 grid remains a popular choice (as it is supported by the native Android OS), it is worth studying whether having more touch dots on a pattern grid leads to stronger security. In this experiment, we first ranked all possible patterns for each grid setting in ascending order according to their complexity scores. We then equally divided the patterns into three groups, simple, medium and complex, and asked our participants to randomly select 20 patterns from each group for evaluation. We report the success rate of our attack within five attempts. In the experiments, we have adapted our algorithms for each grid setting by adjusting the algorithm parameters (such as the line direction numbers).

Figure 26 shows the success rate of our attack for different grids. Similar to a 3×3 grid, our approach achieves a higher success rate for complex patterns over simple ones. On average, we can crack 90% of the complex patterns. We observed that a grid with more dots does provide stronger protection. For complex patterns, the success rate of our attack drops from



(a) pattern-based interface



(b) pin-based interface

Figure 27. Pattern-based and pin-based authentication interfaces from Xiaomi MI4 phone.

95% on a 4×4 grid to 87% on a 6×6 grid. For simple patterns, the success rate of our attack drops from 85% on a 4×4 grid to 75% on a 6×6 grid. This is because a fingertip trajectory can be mapped to a larger number of candidates on a grid with more dots. For instance, the pattern shown in Figure 2 (f) can be mapped to 55 candidate patterns on a 6×6 grid as opposite to 5 on a 3×3 grid. Overall, our attack can crack over 75% (up to 95%) of the patterns within five attempts. One of the purposes of introducing pattern grids with more dots is to allow users to use more complex patterns. However, this experiment suggests that complex patterns remain less security on these grids under our attack.

I. FIXED: Evaluation of Cracking PIN-based Passwords

Result 9: PIN-based passwords are still vulnerable under our attack. This can be proved that we can break over 85% of the passwords within five attempts using a simple variant method based on our approach.

Recall the basic idea of our attack, pattern lock can be reconstructed by tracking the fingertip movement from the target video. Intuitively, our method can also infer to the PIN-base passwords by analyzing the pattern which is generated during unlocking the device. To prove this intuitively imagine, we total collected 30 4-digital PIN-based passwords. For each password, we ask our participants to keystroke the passwords on a Xiaomi MI4 phone. Figure 27 shows the pattern-based and pin-based authentication interface of MI4 phone. We use an iPhone4S phone to record the video respectively on the left-front, front and right-front of the target device from a distance of 2 meters.

Variant Methodology There exists some differences between PIN-based passwords and pattern lock. These are summarized as follow: (1) PIN-base authentication interface has ten dots that is different from pattern lock, making the layout of interface of the two mechanisms different; (2) the dots on PIN-based interface can be typed repeated while they can only be visited once on pattern-base interface. The above two differentiations make the structure of the pattern generated from PIN-based interface more complicated than the one from pattern-based interface. Our preliminary experiments suggest that we can reconstruct the trajectory of PIN-based password by connecting the touching points⁶. We can get the location

⁶Touching points are those points that tracked when fingertip touches the screen.

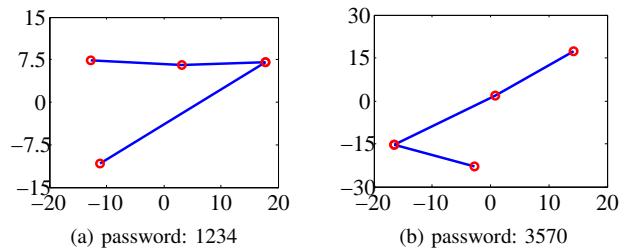


Figure 28. Examples of tracked fingertip trajectory comprised of touching points. The red circle represents the touching points.

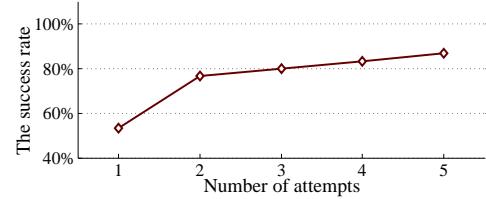


Figure 29. Success rate of cracking PIN-based passwords with different number of attempts.

of touching point by the up-and-down motion direction of the fingertip inspired by the solution in [?]. To reconstruct PIN-based passwords, we increased the geometric information including both direction and length information to adjust the variant attack method. Figure 28 shows the reconstructed fingertip trajectory using the variant method.

Results Figure 29 shows the performance on cracking 4 digital passwords within different number of attempts. With first attempt, the success rate is only 53.3% as most of tested passwords have two or more candidate patterns. As the increase of trial number, the success rate is significantly increased, and they respectively are 76.7%, 80%, 83.3%, 86.7% within 2, 3, 4, 5 attempts. Within five attempts, our method failed to crack four passwords. The reason that we failed two passwords is because they have more than five candidate patterns and for remaining two ones is that we detected the wrong location of one touching point due to the blue motion of the video footage. We believe that the expert attacker is able to improve the success rate using more attempts.

VII. DISCUSSIONS

A. FIXED: Countermeasures

The success of our attack depends on three factors: (1) knowledge of the pattern grid; (2) a decent quality video footage allowing the algorithm to track the fingertip movement; (3) successfully identifying a video segment that captures the entire process of pattern drawing.

Potential countermeasures For the first factor, the attacker can obtain relevant information via analyzing a device installed with the same operating system and applications as the target. Randomization techniques such as randomized pictures [? ?] could be a solution for the first factor. However, randomization-based solutions often come at the cost of poorer usability. This can prevent them to be used at a large scale. Regarding the second factor, there are ways, such

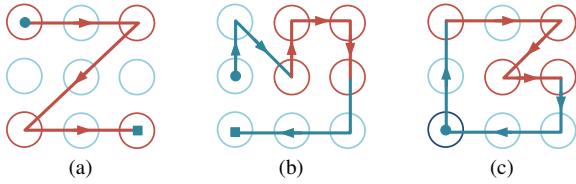


Figure 30. Examples of the variant pattern lock mechanism for defending against the video-based attacks. The red patterns are the correct pattern drawn when enrolment, and the blue dot in (c) shows that the dot is visited twice. (a) and (c) show that a previously unvisited dot can be bypassed if it is part of a horizontal, vertical or diagonal line segment of the pattern.

as KALEIDO [?], to prevent unauthorized videotaping by dynamically changing the colour and brightness of the screen to confuse the filming camera. A non-technical solution for this aspect would be to educate users to fully cover their fingers when drawing a pattern. But doing this on a large-screen device could be awkward especially when the device is held by one hand. For the third factor, the attacker’s solution depends on the type of the pattern. For a screen lock, pattern drawing is the first activity (except for receiving a phone call or making an emergency call) when the device is retrieved. Therefore, identifying the video segment is straightforward. When the pattern is used by applications, we have observed that users typically pause for a few seconds before or after entering the pattern. Therefore, an experienced attacker should also be able to identify the video segment in case our automatic algorithm (presented in Section IV-A) fails to do so. A potential countermeasure is to mix pattern unlocking with other on-screen activities. For examples, before and after pattern drawing, the system can ask the user to type in a sentence using a Swype-like method or to draw some graphical shapes. The problem of this approach is it may annoy users by asking them to do more, especially for screen unlocking – an activity that is performed many times a day.

A Specific Case To invalid this attack, we propose and implement a prototype of a specific countermeasures by disable the third factor presented in last paragraph. The basic idea is that before or after drawing the correct pattern, user can draw any pattern for the purpose of confusing the attacker. In other words, if the pattern drawn by user contains the correct pattern, the system would authenticate successfully. To this end, we redesign the rules for creating a pattern to ensure that each dot on the pattern grid can be visited repeatedly. Further, for purpose of increasing the number of candidate patterns, we also change the rule that a previously unvisited dot can be bypassed if it is part of a horizontal, vertical or diagonal line segment of the pattern.

Figure 30 shows some patterns drawn using our redesign method. From this we can see that the redesign keep the usability of original pattern-based authentication scheme as maximum as possible. Moreover, comparing to original system, this can frustrate shoulder-surfing attack as well as video-based attack as the attacker cannot recognize the correct pattern in spite of see or reconstruct the pattern. Furthermore, this improved scheme precede biometric-based scheme such as face recognition [?] and iris recognition [?] in terms of

revocability.

B. Implications

While pattern lock is preferable by many users [?], this work shows that it is vulnerable under video-based attacks. Our attack is able to break most patterns in five attempts. Considering Android allows five failed attempts before automatically locking the device, our work shows that this default threshold is unsafe. We demonstrated that, in contrast to many users’ perception, complex patterns actually do not provide stronger protection over simple patterns under our attack.

It is worth mentioning that our approach is only one of the many attacking methods that researchers have demonstrated. Examples of these attacks include video-based attacks on keystroke-based authentication [? ?], sensor-based attacks for pattern lock [?]. Authentication methods that combine different authentication methods [? ? ?] to constantly checks the user’s identity could be a solution.

VIII. RELATED WORK

Our work lies at the intersection between computer vision based attacks and cracking graphical- and touch-based authentication methods. This work brings together techniques developed in the domain of computer vision and motion tracking to develop a new attack.

Computer Vision-based Attacks No work has targeted using video footage to crack Android pattern lock and this is the first to do so. Our work is inspired by the work presented by Shukla *et al.* [?] on video-based attacks of PIN-based passwords. In addition to addressing the new challenges highlighted in Section I, our work differs to their approach in two ways. Firstly, we target a different authentication method, i.e. graphical-based passwords are fundamentally different from PIN-based passwords. Secondly, our approach does not require knowledge of the size of the screen or the grid. Other work in the area including [?] which attacks PIN-based passwords by analyzing how the screen brightness changes when entering a password. But the subtle changes of the screen brightness can be dramatically affected by the lighting condition. This restricts the application of their approach. There is a body of work using reflections to recover information typed by the user [? ? ? ?]. These schemes require having a clear vision of the content displayed on the screen while our approach does not have such a requirement.

Cracking Graphical-based Passwords Aviv *et al.* demonstrated that it is possible to reconstruct a locking pattern by analyzing the oily residues left on the screen [?]. This method is highly restricted as oily residues can be messed up by any on-screen activities after pattern drawing. Zhang *et al.* exploit the WiFi signal interferences caused by finger motions to recover patterns [?]. Their method requires a complex setup and is highly sensitive to moving objects of the environment because the WiFi signal can be disrupted by a moving object.

Attacks on Touch-based Authentication Ballard *et al.* implemented a forgery attack on handwriting authentication [?]. Using a small number of training examples, they achieve a high success rate for this attack. More recently, Serwadda

et al. show that a simple robot can achieve high penetration rates against touch-based authentication systems by analyzing on-screen gestures including swiping and zooming [?]. In this paper, we present a new, video-based attack for graphical-based passwords. Research in this area all demonstrates the need for a closer look of the security risks of touch-based authentication.

Study of Android Pattern Lock Uellebenk *et al.* study how people use Android pattern lock on a daily basis [?]. They found that although there is a large number of Android patterns, in practice many people only use a small set of them due to users' bias in generating patterns. Løge explored the correlation between human's characteristics (e.g. ages and genders) and the choice of patterns [?]. Her study shows that users have a bias in selecting the starting dot to form a pattern and people tend to use complex patterns for sensitive applications.

Motion Tracking In addition to TLD, there are other methods proposed in the past for tracking object motions. Some of them apply image analysis to track the hand and gesture motions from video footage [? ? ?]. In this paper we do not seek to advance the field of motion tracking. Instead we demonstrate

that a new attack can be built using classical motion tracking algorithms. We show that the attack presented in this work can be a serious threat for Android pattern lock.

IX. CONCLUSIONS

This paper has presented a novel video-based side-channel attack for Android pattern lock. The attack is based on a video filmed a distance of 2 meters away from the target device using a mobile phone camera. The attack is achieved by employing a computer vision algorithm to track the fingertip movement from the video, and then using the geometry information of the fingertip movement trajectory to identify the most likely patterns to be tested on the target device. Our approach was evaluated using 120 unique patterns collected from independent users and some of the most complex patterns. The experimental results show that our approach is able to successfully crack over 95% of the patterns in five attempts. We show that, in contrast to many people's belief, complex pattern actually provides weaker protection over simple patterns under our attack. Our study suggests that Android pattern lock is vulnerable to video-based side-channel attacks.