

# Sloth: Key Stretching and Deniable Encryption using Secure Elements on Smartphones

Daniel Huguenroth  
University of Cambridge  
dh623@cam.ac.uk

Sam Cutler  
The Guardian  
sam.cutler@theguardian.com

Alberto Sonnino  
MystenLabs & University College London  
alberto@mystenlabs.com

Alastair R. Beresford  
University of Cambridge  
arb33@cam.ac.uk

## ABSTRACT

Privacy enhancing technologies must not only protect sensitive data in-transit, but also locally at-rest. For example, anonymity networks hide the sender and/or recipient of a message from network adversaries. However, if a participating device is physically captured, its owner can be pressured to give access to the stored conversations. Therefore, client software should allow the user to plausibly deny the existence of meaningful data. Since biometrics can be collected without consent and server-based authentication leaks metadata, implementations typically rely on memorable passwords for local authentication.

Traditional password-based key stretching lacks a strict time guarantee due to the ease of parallelized password guessing by attackers. This paper introduces Sloth, a key stretching method leveraging the Secure Element (SE) commonly found in modern smartphones to provide a strict rate limit on password guessing. While this would be straightforward with full access to the SE, Android and iOS only provide a very limited API. Sloth utilizes the existing developer SE API and novel cryptographic constructions to build an effective rate-limit for password guessing on recent Android and iOS devices. Our approach ensures robust security even for short, randomly-generated, six-character alpha-numeric passwords against adversaries with virtually unlimited computing resources. Our solution is compatible with approximately 96% of iPhones and 45% of Android phones and Sloth seamlessly integrates without device or OS modifications, making it immediately usable by app developers today. We formally define the security of Sloth and evaluate its performance on various devices.

Finally, we present HiddenSloth, a plausibly-deniable encryption scheme leveraging Sloth. It provides multi-snapshot resistance against adversaries who can covertly capture its on-disk content multiple times.

## KEYWORDS

key stretching, deniable encryption, secure element, password hashing, Android, iOS

## 1 INTRODUCTION

Smartphones store very sensitive information, ranging from corporate secrets to our most intimate messages and pictures. All users rely on their smartphones to protect their data when devices are lost, stolen, or seized. This is accomplished at the device level through full disk encryption. Apps with strong security requirements typically introduce an additional layer of encryption. Examples include password managers and wallets for digital currency which store sensitive data and need to work offline.

While such apps often allow access using biometrics as a convenience, they typically require a master password as a fallback since sensor readings might be unreliable or to facilitate access by different people. In some cases, apps may intentionally avoid biometric authentication due to concerns about false positives or a desire to prevent involuntary unlocking of the secret: a user can deny knowledge of a password, whereas fingerprints can be obtained without their consent.

The whistleblowing system CoverDrop [1] (presented at PETS '22) is an example of a system requiring plausibly-deniable encryption. In particular, the set of potential sources that might have leaked confidential information is necessarily small. This allows adversaries to enumerate suspects in order to then investigate them and their devices. Since CoverDrop is installed as part of a regular news app, its presence on disk is not suspicious per se. Importantly, in the plausibly-deniable encryption scheme the adversary cannot distinguish whether they are given a wrong passphrase or the user has not used the CoverDrop feature. However, the encryption scheme presented in the CoverDrop design has limitations: its strength is proportional to the size of the stored data, it has no formal security model, and it does not work on iOS devices.

Passwords present a significant challenge: user-chosen passphrases are typically low-entropy, something which is particularly true on mobile devices with small on-screen keyboards. This makes apps that use passwords vulnerable to brute-force attacks. As a countermeasure, apps typically use key stretching schemes to increase the cost for an adversary. Such schemes increase the computational costs, e.g. PBKDF2, or use memory-hard functions, e.g. Argon2. Unfortunately, we have to choose high (and thus expensive) parameters in anticipation that the adversary has access to many computers.

Modern smartphones come with a built-in Secure Element (SE), a separate chip that is hardened against physical tampering and side-channel attacks, and contains its dedicated CPU, RAM, storage, and operating system. The SE communicates with the rest of the

system by message passing which means that even a local adversary with kernel-level privileges and hardware access cannot break the security properties of the SE. While the smartphone vendor can deploy custom code to the SE, including to rate limit operations, app developers only have access to a limited API that provides standard cryptographic operations.

Nevertheless, in this paper, we show that it is possible to build an effective key stretching scheme with SEs as found on smartphones today by using the limited bandwidth of the chip as an intentional bottleneck. Under the assumption that the SE is secure, this places an absolute time cost on each attempt to guess a password since the adversary is forced to make all guesses via the SE, and access to additional computational power is of no benefit. This is in contrast to traditional key stretching schemes where difficult trade-offs are required. For example, an increase in the number of rounds of a key stretching scheme increases the cost of a successful attack for a password of a given strength but also increases the time taken by the app to authenticate the legitimate user.

In this paper, we also explore particularly sensitive scenarios where encryption in itself might not be enough, and we consider the case where the user may need to deny the existence of any meaningful content at all. In a deniable encryption (DE) scheme the decryption function always fails unless it is called with the correct password (if any such password exists). However, decryption failure is not distinguishable from trying to decrypt any (random) byte string that does not contain hidden information. This requires an app to initialize on installation any local state using a randomly-encrypted ciphertext. Such DE schemes can be used to build deniable storage containing hidden volumes. DE schemes inherently rely on passwords as a means of deniable authentication and our work provides an effective SE-backed DE scheme. We extend our DE scheme so that it leaks no information even if an attacker is given multiple storage snapshots; e.g., access to several snapshots of the device state might be available through cloud backups.

*We make the following contributions.*

- (1) We design the Sloth key stretching scheme to provide effective rate-limiting against powerful adversaries. Our scheme is practical since it uses existing APIs and does not require software or hardware modifications.
- (2) We evaluate these schemes on different Android and iOS phone models and show that even short passwords of six alpha-numerical characters provide strong security against the most powerful adversaries while keeping user-initiated unlock operations faster than 1 second.
- (3) We formally capture the security of schemes using SEs by introducing the model of a wall-time-bounded (WT) adversary with oracle access to the SE; we prove the Sloth scheme secure in this setting.
- (4) We design the HiddenSloth DE scheme that can withstand multi-snapshot adversaries and prove it secure.

*We had to overcome the following key challenges.* Mobile apps are sandboxed and can only access a limited API. In particular, SEs do not provide explicit rate-limiting APIs and apps cannot run custom code on the SE. Documentation for both platform is limited requiring both an availability survey (Appendix A) and exploratory

engineering. Our design required bespoke formalization to allow for security analysis which is different from the one used for computationally and memory-hard key stretching functions.

## 2 BACKGROUND

Hardware support for executing code in a secure context is available in most modern smartphones. Trusted Execution Environments (TEE) are the first architecture that achieved wide integration. TEE implementations, such as ARM TrustZone [35], run on the main chip, but in a privileged context so that even a compromised Kernel cannot manipulate them. However, TEEs generally come with limited protection against side-channel and physical attacks. Secure Elements (SE) are standalone components with their dedicated CPU, memory, and storage. This provides stronger isolation and protection against physical attacks.

Android allows developers to create and use keys via its API which abstracts a *Keymaster* interface to manage keys. From API level 23 (Android 6.0 M, released 2015) application developers can verify that keys are stored inside secure hardware. For phones released around 2015, this would typically mean that they are managed by a Keymaster implementation inside a TEE. However, from API level 28 (Android 9 P, released 2018), Android supports the StrongBox Keymaster implementation which must be implemented using an SE [20, 22]. Google has supported StrongBox in its flagship models since the Pixel 3 release in 2018. Other devices might have included SEs before this, but app developers would not have been able to use them.

Apple refers to SEs in their devices as a *Secure Enclave* and they have been supported in iOS devices since the iPhone 5S (released 2013) [4]. From the beginning, they have been used to protect biometric data and secure device encryption. However, they only were exposed to app developers with iOS 13 (released 2019) [7].

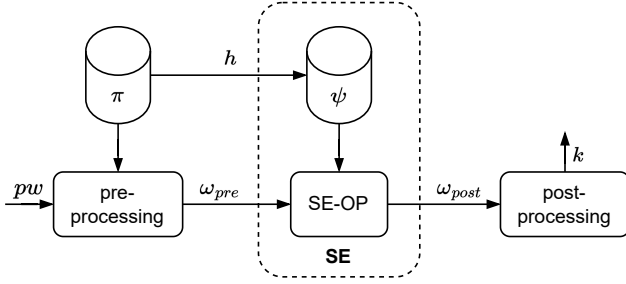
While TEEs and SEs allow for a much smaller implementation and attack surface, they are not impenetrable. Researchers have successfully extracted private keys from TrustZone implementations from Qualcomm [38] and Samsung [39]. Intel’s SGX extension is similar in broad terms to TrustZone where there are many documented attacks [31].

Appendix A provides a survey on the availability of SEs in modern smartphones, their APIs, and their performance. We find that 96% of iPhones and 45% of Android devices provide SE access to app developers. Virtually all recent devices have at least TEE-backed key protection.

### 2.1 Android and iOS APIs

Code that runs on the SE is (by design) executed without any control by the operating system. Therefore, platforms do not allow developers to run custom code. Instead, they offer access to specific cryptographic operations through APIs. These APIs serialize the user request and send it to the SE where it is parsed and executed. The result is returned similarly.

We discuss the API on iOS first as it consists of very few operations. This provides a minimal attack surface, but as we will see later, it also stands in the way of efficient software implementation. At the time of writing, iOS only supports storing private EC keys of type P256 [8]. Once created inside the SE, the API provides methods



**Figure 1: Overview of all Sloth schemes.** The user password,  $pw$ , is first pre-processed on the device with access to device state  $\pi$ , producing  $\omega_{pre}$ ; the SE executes a cryptographic function over  $\omega_{pre}$  using its internal (hidden) state  $\psi$ ; the SE’s output,  $\omega_{post}$ , is then returned to the device; the device can then post-process  $\omega_{post}$ , for example by applying a key derivation function, resulting in  $k$  as the final output.

to sign data, output the corresponding public key, and perform key agreement via ECDH. The API also provides methods for encrypting and decrypting data. However, the documentation is not clear on whether this happens entirely within the SE, or whether the ECDH result is shared with an AES engine outside the chip. We provide additional discussion on why the iOS API does not allow simpler schemes in Appendix B.

The Android API offers more operations and key types to developers. It uses the Keymaster API which abstracts the actual key handling from the user. Device vendors implement the Keymaster HAL which serializes the API calls and communicates with the backend. This backend could be a service running in a TEE or SE. For our paper, we are interested in the StrongBox Keymaster implementation which requires using a SE [22]. It guarantees the availability of the following algorithms: RSA-2048, AES-128/256, ECDSA P-256, ECDH P-256, HMAC-SHA256, and 3DES. The availability of symmetric cryptography allows us to come up with a simpler design for Android devices. However, it also increases the complexity of the implementation that device manufacturers have to provide. Android recently added a new API `setMaxUsageCount` for OS version 12+ (Android S, API 31) that deletes the key after a given number of operations. While it would be convenient, it is implemented at the OS level<sup>1</sup>, as StrongBox does not have sufficient persistent per key storage and thus cannot manage respective counters internally.

### 3 SYSTEM OVERVIEW

Figure 1 provides a high-level overview common to all our Sloth schemes. Sloth distinguishes between two execution spaces: the general device space and the SE. We assume operations performed on the main device may be controlled by the adversary (e.g. using a local kernel exploit) while we assume the SE is a separate piece of hardware and remains secure (see Section 3.1). The user inputs a password,  $pw$ , in the user space. The sloth schemes start by pre-processing  $pw$  using the state  $\pi \in \Pi$  where  $\pi$  is accessible in user

space; the output of this pre-preprocessing step,  $\omega_{pre}$ , is given as input to the SE. The SE maintains its own secure state,  $\psi \in \Psi$ , with cryptographic keys inaccessible to the device. Each key maintained in  $\psi$  is associated with a unique (and public) key handle  $h$ . The SE accesses  $\psi$  and executes a cryptographic operation SE-Op; its output  $\omega_{post}$  is returned to the user space for further post-processing to derive the final key  $k$ . The key  $k$  can then be used for authentication, full disk encryption, or deniable encryption protocols.

#### 3.1 Threat Model

We assume an adversary can physically capture the device. For instance, they find a lost device on the street or stop the user at a border crossing. Before the capture, we assume that the device can be securely and confidentially operated by the user. After the capture, only the SE resists full-access by the adversary.

The adversary aims to determine whether the device contains any user data that has been encrypted with a password and if so recover the password, and therefore the data. We assume passwords are a sequence of characters, including passphrases. Further, we assume that if the adversary can determine that encrypted information is present on the device, they pressure the user with this evidence and obtain the correct password. In other words, a successful scheme must provide plausible deniability and prevent an adversary from distinguishing between encrypted information and random data.

There are many proposals [18, 26, 29, 30, 34, 36, 41] for plausibly-deniable encryption (DE). They typically feature a multi-volume model where one passphrase unlocks a benign cover volume and another passphrase unlocks a hidden volume with sensitive data. For simplicity, our definitions are single-volume—challenging the adversary to decide between encryption of padded data and a random bitstring of the same pre-defined length. Nevertheless, our simpler model can be exchanged for a multi-volume one from existing work and then wrapped with HiddenSloth to provide multi-snapshot guarantees. Stenographic approaches, i.e. an attacker cannot determine whether any DE was used, typically require subtle interactions with the underlying layers. Therefore, they are not easily compatible with our HiddenSloth scheme. However, both multi-volume models and stenographic techniques allow for our key stretching schemes to be used as drop-in replacements when processing the user passphrase.

User-chosen passwords are typically of low-entropy as users reuse the same passwords across different services and choose easy to guess words and numbers. Therefore, we designed Sloth to allow for short, memorable passphrases that are generated for the user. Since these are randomly generated, we can assume that they are uniformly distributed. In particular, we assume for the remainder of the paper that the probability distribution of passphrases  $\mathcal{X}$  has Shannon entropy  $m$ . I.e. for all passphrases  $pw$  in the support of  $\mathcal{X}$  and for random variable  $X$  drawn according to  $\mathcal{X}$ ,  $\text{PROB}(X = a) = 2^{-m}$ . However, developers integrating with Sloth might also allow user-chosen passphrases and follow our arguments using a *min-entropy* distribution similar as in the HKDF paper [28].

The adversary can also use the secure element SE as a black box and can perform any operations via the available API (through Oracle queries which we denote  $\mathcal{O}_{SE}$ ). However, the adversary cannot

<sup>1</sup><https://cs.android.com/android/platform/superproject/+/master:system/keymint/common/src/tag.rs;l=314-333;drc=ed657df7c7b329fb3d26eb0ce88af92594245ae8>

extract information about the key material from the SE and they cannot clone the SE. Similar guarantees are given by Apple for their Secure Enclave [4] and Google for StrongBox implementations [22]. These claims are taken seriously by vendors: Apple and Google award up to \$250,000 [6] and \$1,000,000 [21], respectively, for the discovery of relevant vulnerabilities.

We introduce an abstract definition of an SE. The SE operates on a state  $\psi$  that can only be (meaningfully) accessed by it. In practice, many SEs have limited internal storage and store an authenticated ciphertext of their state on storage managed by the OS. This encryption is performed using a secret internal key. As the encrypted key blobs are stored outside the SE, this can allow an attacker to perform roll-back attacks which we discuss as a practical limitation in Section 5.3. For our threat model we assume that the adversary can capture all application data, such as code and stored data, as well as the encrypted key blobs. A multi-snapshot adversary might capture the application data multiple times (e.g. through automated cloud backups), but the encrypted key blobs are only available once they physically capture the device.

In our definitions, the explicit passing of the state models this behavior. Definition 1 introduces a generic SE as there exist SEs with different capabilities. Later sections introduce more specialized definitions for SEs with various capabilities.

**Definition 1** (Secure Element). A Secure Element SE operates on a hidden state  $\psi \in \Psi$  using the algorithm  $\text{SE.INIT}$  and possible extension algorithms. It also has a timing function  $T_{\text{ALG}}$  that maps each algorithm  $\text{ALG}$  and its arguments to a wall time cost.  $\text{SE.INIT} : \emptyset \rightarrow \Psi$  initializes an empty state  $\Psi$  with  $T_{\text{SE.INIT}}() = 1$ . Only the SE that initialized the state can operate on it with any of the extension algorithms.

We assume the adversary has a limited *wall time budget*  $B$  that is spent when performing oracle operations. The costs for each operation  $\text{OP}$  is defined by  $T_{\text{OP}}$  and deduced from  $B$  before it is executed. We require that the adversary ends the experiment with a non-negative time budget  $B \geq 0$ . We define for any algorithm (including adversary and challenger):

**Definition 2** (Wall Time Algorithm). A wall time (WT) algorithm is a PPT algorithm  $A$  with an initial wall time budget of  $B_A \in \mathbb{N}$ . During its execution it can perform the operations  $\text{OP}_i(p_1, p_2, \dots)$  using the oracle  $\mathcal{O}_{\text{SE}}$  given  $\sum_i T_{\text{OP}_i}(p_1, p_2, \dots) \leq B_A$ .

In comparison to standard security definitions, the strength is no longer supported by just asymptotic computational effort in  $\text{poly}(\lambda)$ , but also by a concrete wall time budget  $B$ . This is a strong property, as it is independent of advances in processing power or utilizing more machines.

### 3.2 Notations

Throughout the paper, let  $\lambda$  be a freely chosen but fixed security parameter. We write  $x.y \leftarrow z$  to denote the assignment of value  $z$  to the field  $y$  of a named-tuple  $x$ . Table 1 summarizes the symbols used in the rest of the paper.

## 4 THE SLOTH KEY STRETCHING SCHEME

We now describe the design of two concrete variants of our key stretching scheme Sloth. In both cases we leverage the limited

Algorithms and schemes	
HKDF	Key derivation function
$\Xi$	Key stretching scheme
$\Delta$	Deniable encryption scheme
Parameters	
$\lambda$	Security parameter for computational security
$l$	Length of $\omega_{\text{pre}}$ in LongSloth
$n$	Count of $\omega_{\text{pres}}$ in RainbowSloth
Variables	
$\text{pw} \in \mathcal{P}$	The user password (space)
$m$	The Shannon entropy of the password distribution
$\pi \in \Pi$	Storage state (space)
$\psi \in \Psi$	State (space) inside the SE
$h \in H$	Key handle (space) for the SE
$\omega$	An intermediate secret
$k$	The final derived secret key

**Table 1: We use these symbols for algorithms, parameters, and variables in our descriptions throughout the paper.**

throughput of the SE to effectively rate-limit password guessing. The variant LongSloth (Section 4.1) only uses symmetric operations and is simpler. In particular, it uses an HMAC operation and by increasing the required length of its input, we can reduce the guess rate. However, LongSloth is incompatible with the iOS API which does not offer symmetric cryptographic operations. The variant RainbowSloth (Section 4.2) is compatible with both Android and iOS APIs as it only requires support for a DH key exchange, but at the cost of higher complexity and additional assumptions. In RainbowSloth, we make the final key dependent on the result of multiple ECDH operations and thus we can reduce the guess rate by increasing the required number of key agreements.

### 4.1 LongSloth: Simple Key Stretching Scheme

LongSloth exclusively relies on symmetric operations both outside and inside the secure element SE. This enables extremely efficient implementations on Android, but cannot be implemented on iOS (see Section 4.2 for a variant compatible with iOS). Particularly, LongSloth operates on an SE equipped with HMAC support:

**Definition 3** (SE with HMAC Support). A Secure Element with HMAC support  $\text{SE-WITH-HMAC}$  is an SE that has two extension algorithms:

- $\text{SE.HMACKEYGEN} : \Psi \times H \rightarrow \Psi$  generates a new secret key  $k$  and updates  $\psi \in \Psi$  under handle  $h \in H$ :  $\psi.h \leftarrow k$  requiring time  $T_{\text{SE.HMACKEYGEN}}(\psi, h) = \Theta(1)$ .
- $\text{SE.HMAC} : \Psi \times H \times M \rightarrow \{0, 1\}^\lambda$  takes a message  $m \in M$  and a key handle  $h \in H$  and outputs the  $\text{HMAC}(\psi, h, m)$  result of a message  $m \in M = \{0, 1\}^*$  requiring time  $T_{\text{SE.HMAC}}(\psi, h, m) = c_{\text{HMAC}} \cdot |m| = \Theta(|m|)$ , where  $c_{\text{HMAC}}$  is a device specific constant.

**Algorithm 1** The LongSloth protocol with the security parameters  $l, \lambda$  freely chosen, but fixed.

---

```

1: procedure LONGSLOTH.KEYGEN( $\psi, pw, h$ )
2:    $\pi \leftarrow \{\}$ 
3:    $\pi.h \leftarrow h$ 
4:    $\pi.salt \xleftarrow{\$} \{0,1\}^\lambda$ 
5:    $\psi \leftarrow \text{SE.HMACKEYGEN}(\psi, h)$ 
6:    $k \leftarrow \text{LONGSLOTH.DERIVE}(\pi, \psi, pw)$ 
7:   return ( $\pi, \psi, k$ )
8:
9: procedure LONGSLOTH.DERIVE( $\pi, \psi, pw$ )
10:   $\omega_{pre} \leftarrow \text{PWHASH}(\pi.salt, pw, l)$ 
11:   $\omega_{post} \leftarrow \text{SE.HMAC}(\psi, \pi.h, \omega_{pre})$ 
12:   $k \leftarrow \text{HKDF}(\omega_{post})$ 
13:  return  $k$ 

```

---

Algorithm 1 describes the main operations of LongSloth. The procedure LONGSLOTH.DERIVE (Line 9) outputs a potentially existing key. During the pre-processing step (see Figure 1), LongSloth expands a user password to a bit string  $\omega_{pre}$ . This is achieved using the PWHASH operation that reads a salt value from the local storage and then hashes the user password to a variable length output. The output length  $l = |\omega_{pre}|$  is a configurable parameter to limit the guess rate based on the SE throughput rate. Next,  $\omega_{pre}$  is computed by the SE-OP securely inside the SE with a key selected by the key handle  $h$ . This operation is an HMAC producing  $\omega_{post}$  that is then hashed into the final key output  $k = \text{HKDF}(\omega_{post})$  using the hash-based key derivation function HKDF [28].

The generation of a new key  $k$  is done similarly. The procedure LONGSLOTH.KEYGEN (Line 1) first initializes a new state  $\pi$  with the key handle  $h$  and a fresh random salt value. The SE then generates a new HMAC key under  $h$ . It returns the updated state and the key  $k$  by calling LONGSLOTH.DERIVE.

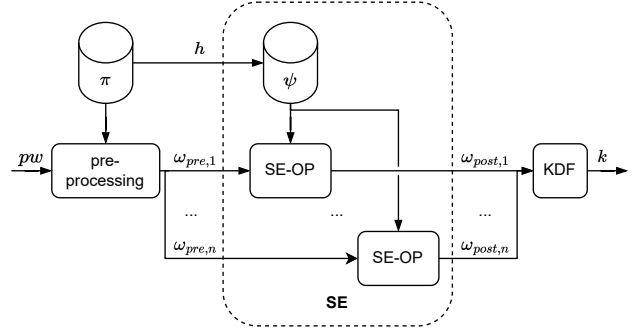
**Security intuition.** LongSloth is designed to withstand an adversary with a wall time budget  $B < 2^m \cdot l \cdot c_{\text{HMAC}}$ . That is, the adversary would need to call  $2^m$  times the operation SE.HMAC( $x$ ) with  $|x| = l$  (see Definition 2) to try all possible user passwords. Intuitively, the security of LongSloth relies on the observation that such an adversary is unable to distinguish a key  $k$  generated by LongSloth from a random bit string. Section 6 provides a formal security analysis based on the observation that HMAC is a PRF [12].

## 4.2 RainbowSloth: Key Stretching for iOS

RainbowSloth relies on an SE equipped with support for Elliptic Curve Diffie-Hellman (ECDH) key exchanges (Definition 4). RainbowSloth is more complex than LongSloth but is compatible with both Android and iOS.

**Definition 4** (SE with ECDH Support). A Secure Element with Elliptic Curve Diffie-Hellman support SE-WITH-ECDH is a SE that has two extension algorithms:

SE.ECDHPrivKeyGen :  $\Psi \times H \rightarrow \Psi$  generates a private EC key  $k$  and updates  $\psi \in \Psi$  under handle  $h \in H$ :  $\psi.h \leftarrow k$  requiring time  $T_{\text{SE.ECDHPrivKeyGen}}(\psi, h) = \Theta(1)$ .



**Figure 2: RainbowSloth scheme with a sequence of  $\omega_{pre,i}$  inputs. The dashed area indicates the SE. Since there is no parallelism, all SE-Op are executed sequentially.**

$\text{SE.ECDH} : \Psi \times H \times \text{pub} \rightarrow \{0,1\}^\lambda$  takes a EC public key  $\text{pub}$  and a key handle  $h \in H$  and outputs the key exchange result  $\text{ECDH}(\psi, h, \text{pub})$  requiring time  $T_{\text{SE.ECDH}}(\psi, h, \text{pub}) = c_{\text{ECDH}} = \Theta(1)$ , where  $c_{\text{ECDH}}$  is a device specific constant.

In contrast to LongSloth, RainbowSloth computes a sequence of fixed-sized public keys as its pre-secrets  $\omega_{pre,i}$  (see Figure 2). The processing speed of each ECDH operation is constant since the public key size is defined by the underlying elliptic curve. Instead of increasing the input length of each operation, RainbowSloth increases the required number of SE-Op executions to achieve a given throughput limit.

Algorithm 2 describes the main operations of RainbowSloth. The procedure RAINBOWSLOTH.DERIVE decrypts a potentially existing key. It first expands the user password into multiple bit strings  $\omega_{pre,i}$  with  $i \in [1, n]$  calling PWHASH. Those bit strings are then formatted into P-256 public keys (Line 13). The SE uses the key handle  $h$  to select a private P-256 key and computes an ECDH operation with each  $\omega_{pre,i}$  (Line 14). Since the SE has no parallelism, these operations take place one after another. The resulting values  $\omega_{post,i}$ ,  $i \in [1, n]$  are then jointly hashed into a final key  $k$  such that each contributes to all bits of  $k$ .

The generation of a new key  $k$  is done similarly. The procedure RAINBOWSLOTH.KEYGEN (Line 1) first initializes a new state  $\pi$  with the key handle  $h$  and a fresh random salt value. The SE then creates a new P-256 key under  $h$ . It returns all updated states and the key  $k$  by calling RAINBOWSLOTH.DERIVE as described above.

**Security intuition.** RainbowSloth is designed to withstand an adversary with wall time budget  $B < 2^m \cdot n \cdot c_{\text{ECDH}}$  (see Definition 2); the adversary would need to call  $2^m \cdot n$  times SE.ECDH to try all possible user passwords. Similarly to LongSloth, the security of RainbowSloth relies on the observation that the adversary is unable to distinguish a key  $k$  generated by RainbowSloth from a random bit string. Section 6 provides a formal security analysis assuming the SE runs the ECDH operations in a group where the generalized decisional Diffie-Hellman problem is hard [11].

**Hashing into the P-256 public key space.** RainbowSloth requires an operation HASHToP256 that maps any bit-string  $\text{seed} \in$

**Algorithm 2** The RAINBOWSLOTH protocol with the security parameters  $n, \lambda$  freely chosen, but fixed.

---

```

1: procedure RAINBOWSLOTH.KEYGEN( $\psi, pw, h$ )
2:    $\pi \leftarrow \{\}$ 
3:    $\pi.h \leftarrow h$ 
4:    $\pi.salt \xleftarrow{\$} \{0, 1\}^\lambda$ 
5:    $\psi \leftarrow \text{SE.ECDHPRIVKEYGEN}(\psi, h)$ 
6:    $k \leftarrow \text{RAINBOWSLOTH.DERIVE}(\pi, \psi, pw)$ 
7:   return  $(\pi, \psi, k)$ 
8:
9: procedure RAINBOWSLOTH.DERIVE( $\pi, \psi, pw$ )
10:   $l \leftarrow n \cdot (\frac{256}{8})$ 
11:   $\omega_{pre,1} \parallel \dots \parallel \omega_{pre,n} \leftarrow \text{PwHASH}(\pi.salt, pw, l)$ 
12:  for  $i = 1 \dots n$  do
13:     $x \leftarrow \text{REHASHTOP256}(\omega_{pre,i})$ 
14:     $\omega_{post,i} \leftarrow \text{SE.ECDH}(\psi, \pi.h, x)$ 
15:   $k \leftarrow \text{HKDF}(\omega_{post,1} \parallel \dots \parallel \omega_{post,n})$ 
16:  return  $k$ 

```

---

$\{0, 1\}^*$  to a valid P-256 public key. Each P-256 public key can be represented by 256 bits for its X-coordinate and a bit that determines the Y-coordinate [17]. However, not all possible 257-bit strings refer to valid public keys, as the key space only has size  $2^{256}$ . We designed the REHASHTOP256 algorithm and use it in our practical evaluation (Section 7).

The REHASHTOP256 algorithm (see Algorithm 3) considers the SEC-1 octal representation of P-256 curve points with point compression [17]. The representation for P-256 keys starts with a byte that is either  $0x02$  or  $0x03$  for the Y-coordinate information bit followed by 32 bytes for the X-coordinate. The algorithm repeatedly hashes the seed string and a counter to a candidate octet array. The first octet is then set as  $0x02$  or  $0x03$  based on the parity of the original value of the first octet. We try to convert each candidate array to a P-256 public key using the SEC1OCTETIMPORTP256 algorithm [17, 2.3.4]. If it returns “invalid” ( $\perp$ ) or the point at infinity ( $\mathcal{O}$ ), the counter is incremented and the algorithm is repeated. Otherwise, a valid P-256 public key is returned. A similar approach is sketched in the Elligator paper [13, 1.4].

The runtime of this algorithm depends on finding a valid public key representation in one of its iterations. Since the output of the KDF can be considered pseudo-random, each iteration is independent and roughly half of the candidate octet arrays are invalid. We thus expect that the algorithm terminates with less than 10 iterations in  $1 - 0.5^{10} > 99.9\%$  of all cases. Hashing with a counter instead of re-hashing the seed avoids falling into small loops that consist only of “bad seeds”. While long executions are unlikely, we are unaware of proofs of this algorithm’s polynomial-time termination.

Alternatively, one might build this operation using the Elligator Squared ( $E^2$ ) [42] technique.  $E^2$  maps an elliptic curve point to a bit string that is indistinguishable from random. Such a bit string can then be mapped back to an elliptic curve point. While it can be computed efficiently [9],  $E^2$  bit strings require more space. A more practical concern is that we are not aware of any freely available implementation. However, if deterministic runtime is required,  $E^2$  can be used as a drop-in replacement for REHASHTOP256.

**Algorithm 3** The REHASHTOP256 algorithm that maps any bit-string  $seed \in \{0, 1\}^*$  to a valid P-256 public key.

---

```

1:  $counter \leftarrow 0$ 
2: while true do
3:    $arr \leftarrow \text{HKDF}(seed \parallel counter)[0 : 32]$ 
4:    $arr[0] = 0x02 \mid (arr[0] \& 0x01)$ 
5:    $x = \text{SEC1OCTETIMPORTP256}(arr)$ 
6:   if  $x \neq \perp$  and  $x \neq \mathcal{O}$  then return  $x$ 
7:    $counter \leftarrow counter + 1$ 

```

---

## 5 THE SLOTH DENIABLE ENCRYPTION SCHEME

Our generic methods to store and derive a key  $k$  can be used to build a deniable encryption (DE) scheme. We require a DE scheme to have three methods: (i) an INIT method that initializes a new storage of a given maximum size, (ii) an ENCRYPT method that stores data encrypted with a given password, and (iii) a DECRYPT method that retrieves and decrypts data from storage if there is any saved with the given password (and otherwise fails).

We present two variants of our deniable encryption scheme, 1S-HiddenSloth (Section 5.1) and MS-HiddenSloth (Section 5.2). 1S-HiddenSloth always pads the (potentially empty) payload to a fixed sized and then encrypts it with either a derived passphrase or a random key. Thus, it withstands an adversary capable of capturing a single snapshot of the user’s device. MS-HiddenSloth wraps 1S-HiddenSloth into another layer of symmetric encryption using a key stored in the SE that is frequently rotated. This technique allows re-encryption without user interaction resulting in different ciphertexts that mimic encryption of new content by the user. Thus, MS-HiddenSloth withstands an adversary capable of capturing multiple snapshots at different points in time.

### 5.1 1S-HiddenSloth: Single-Snapshot

For most real-world scenarios, the single-snapshot adversary is the most realistic threat model. An adversary gaining access to a phone through finding or stealing it has no prior snapshot or knowledge of the phone state, and hence cannot perform multi-snapshot attacks. Similarly, an adversary who confiscates a device can rarely do so covertly. If a user suspects an adversary may have taken a snapshot of their phone, they can reset the phone and thus revert to a single-snapshot scenario. Nevertheless, there may be situations where an adversary can obtain multiple copies of the phone state over time; this is addressed in Section 5.2.

Algorithm 4 presents 1S-HiddenSloth, a single-snapshot resistant deniable encryption scheme. When calling the INIT procedure, the algorithm allocates a storage file *blob* and fills it with random bytes to the given size (plus an allowance for overhead). This is implemented by encrypting a zeroed payload using a randomly chosen password (Line 5). Upon encrypting new user data by calling the ENCRYPT procedure, the user provides a passphrase  $pw$  used to derive a cryptographic key  $k$  (using either LongSloth or RainbowSloth, see Section 4). The payload is then prefixed with a 4-bytes integer representing its length and padded to the maximum size  $s$ . The resulting byte string is encrypted using the key  $k$ , for instance using AES-GCM. For the DECRYPT operation, the key  $k$

**Algorithm 4** The *deniable encryption against single-snapshot* DESS protocol (1S-HIDDEN SLOTH) for maximum data size  $s \leq 2^{31}$  and security parameter  $\lambda$  freely chosen, but fixed. The underlying SLOTH protocol can be instantiated with either variant.

---

```

1: procedure DESS.INIT( $\psi, h, s$ )
2:    $pw \xleftarrow{\$} \{0, 1\}^\lambda$ 
3:    $\psi, \pi, k \leftarrow \text{SLOTH.KEYGEN}(\psi, pw, h)$ 
4:    $\pi.s \leftarrow s$ 
5:    $\pi \leftarrow \text{DESS.ENCRYPT}(\pi, \psi, pw, [ ])$ 
6:   return  $(\pi, \psi)$ 
7:
8: procedure DESS.ENCRYPT( $\pi, \psi, pw, data$ )
9:    $k \leftarrow \text{SLOTH.DERIVE}(\pi, \psi, pw)$ 
10:   $x \leftarrow \text{UINT32}(|data|) \parallel data \parallel 0^{(|\pi.s| - |data| - 4)}$ 
11:   $\pi.iv \xleftarrow{\$} \{0, 1\}^\lambda$ 
12:   $\pi.blob, \pi.tag \leftarrow \text{AE.ENC}(k, iv, x)$ 
13:  return  $\pi$ 
14:
15: procedure DESS.DECRYPT( $\pi, \psi, pw$ )
16:   $k \leftarrow \text{SLOTH.DERIVE}(\pi, \psi, pw)$ 
17:   $x \leftarrow \text{AE.DEC}(k, \pi.iv, \pi.blob, \pi.tag)$ 
18:  if  $x = \perp$  then return  $\perp$ 
19:   $s' \leftarrow \text{UINT32}(x[4 : 4])$ 
20:  return  $x[4 : 4 + s']$ 

```

---

is derived analogously and the algorithms attempt to decrypt the beginning of the file. If this operation fails, an adversary will not be able to tell whether the passphrase was wrong or there was no previous call to ENCRYPT at all.

## 5.2 MS-HiddenSloth: Going Multi-Snapshot

Sometimes the storage state is available to an adversary on multiple occasions. This may happen if an app's data is backed up to the cloud to protect against device loss. Intuitively, 1S-HiddenSloth does not withstand such an adversary because changes in the ciphertext leak whether the storage has been overwritten between the device capture events.

MS-HiddenSloth overcomes this limitation and allows the adversary to capture the storage state  $\pi$  multiple times. However, our model restricts the adversary to only access the SE once they finally gain physical access to the device. That is, the adversary can perform SE calls only during the last capture event. We believe this restriction is practical as smartphones are usually with the user and hence not likely available for covert access by the adversary<sup>2</sup>.

MS-HiddenSloth requires a SE that supports symmetric encryption. This primitive is available directly on Android. On iOS it can be emulated by performing ECDH with a known public key and then use the resulting shared secret for AES-GCM<sup>3</sup>.

<sup>2</sup>Also, if the adversary had such local privileged access multiple times they could simply install malware that records the keyboard and screen.

<sup>3</sup>This is conveniently provided by the API as the `SecureKeyAlgorithm` variant `eciesEncryptionCofactorVariableIVX963SHA256AESGCM`.

**Definition 5** (SE with Symmetric Encryption). A Secure Element with symmetric encryption support SE-WITH-SYMMENC is an SE that has three extension algorithms:

SE.SYMMKEYGEN :  $\Psi \times H \rightarrow \Psi$  generates a new secret key  $k$  and updates  $\psi \in \Psi$  under handle  $h \in H$ :  $\psi.h \leftarrow k$  requiring time  $T_{\text{SE.SYMMKEYGEN}}(\psi, h) = \Theta(1)$ .  
SE.SYMMENC :  $\Psi \times H \times IV \times M \rightarrow C$  takes an initialization vector  $iv \in IV$  and a key handle  $h \in H$  and outputs the ciphertext  $\text{SYMMENC.ENCRYPT}(\psi.h, iv, m) = c \in C$  for the message  $m \in M$  requiring time  $T_{\text{SE.SYMMENC}}(\psi, h, m) = \Theta(|m|)$ .  
SE.SYMMDEC :  $\Psi \times H \times IV \times C \rightarrow M$  takes an initialization vector  $iv \in IV$  and a key handle  $h \in H$  and outputs the message  $\text{SYMMENC.DECRYPT}(\psi.h, iv, c) = m \in M$  for the ciphertext  $c \in C$  requiring time  $T_{\text{SE.SYMMENC}}(\psi, h, c) = \Theta(|c|)$ .

To protect against a multi-snapshot adversary MS-HiddenSloth wraps the storage in another layer of encryption guarded by an additional symmetric key  $k_{SE}$  held in the SE's secure state under handle  $h'$ . The outer layer is periodically re-encrypted using a new temporary key  $tk$ , which in turn is re-encrypted with a fresh  $k_{SE}$ . The re-encryption process (DEMS.RATCHET) should happen at least as often as the adversary has the opportunity to access the stored data. For instance, for a backed-up app, this would happen after every upload operation. In other cases, every app restart could be used as a trigger event. An important design feature of MS-HiddenSloth is that it does not require the user's password to execute the re-encryption process; that is, the procedure DEMS.RATCHET can be executed entirely in the background without user interaction.

**MS-HiddenSloth algorithm.** Algorithm 5 presents all the operations of MS-HiddenSloth. The initial state is initialized similarly to 1S-HiddenSloth but with the creation of an additional SE secret key ( $k_{SE}$ ) associated with the handle  $h'$  (Line 4). To encrypt new data, MS-HiddenSloth first encrypts them as in 1S-HiddenSloth; it then generates an ephemeral key  $tk$  that is used to re-encrypt the data (Line 12). This ephemeral key is finally encrypted using the SE's secret key generated at Line 4 and stored in the user space  $\pi$ . To decrypt existing data, MS-HiddenSloth first decrypts the ephemeral key  $tk$  (Line 21); it then uses that key to decrypt the outer encryption layer (Line 22) and finally decrypt the data (Line 23). The indirection via  $\pi.tk$  allows to leverage the faster AES engine on the main CPU without sacrificing security.

The re-encryption process (ratchet) works similarly to decryption followed by an encryption operation. It firsts retrieves the temporary key  $tk$  and uses it to decrypt the outer encryption layer. It then generates a new temporary key that is used to re-encrypt the data and thus create a new outer encryption layer; the new temporary key is then encrypted by the SE and persisted in the user space. For performance reasons, MS-HiddenSloth does not perform the re-encryption process over the actual data  $\pi.blob$  inside the SE; it instead re-encrypts the data with a temporary key  $\pi.tk$ . This temporary key  $\pi.tk$  and its related fields  $\pi.tiv, \pi.ttag$  are freshly generated at every encryption (Line 12) and ratchet step (Line 29); they are only stored encrypted with the SE's secret key generated at Line 4.

**Algorithm 5** The *deniable encryption against multi-snapshot* DEMS protocol (HIDDEN-SLOTH) for max. size  $s \leq 2^{31}$  and security parameter  $\lambda$  freely chosen, but fixed. The underlying SLOTH protocol can be instantiated with either variant.

---

```

1: procedure DEMS.INIT( $\psi, h, s$ )
2:    $\pi, \psi, k \leftarrow \text{DESS.INIT}(\psi, h \parallel 0, s)$ 
3:    $\pi.h' \leftarrow h \parallel 1$ 
4:    $\psi \leftarrow \text{SE.SYMMKEYGEN}(\psi, \pi.h')$ 
5:    $\pi.seiv \xleftarrow{\$} \{0, 1\}^\lambda$ 
6:    $\pi.pw \xleftarrow{\$} \{0, 1\}^\lambda$ 
7:    $\pi \leftarrow \text{DEMS.ENCRYPT}(\pi, \psi, \pi.pw, [])$ 
8:   return  $(\pi, \psi)$ 
9:
10: procedure DEMS.ENCRYPT( $\pi, \psi, \pi.pw, data$ )
11:    $\pi \leftarrow \text{DESS.ENCRYPT}(\pi, \psi, \pi.pw, data)$ 
12:    $\pi.tk \leftarrow \text{AE.KEYGEN}(); \pi.tiv \xleftarrow{\$} \{0, 1\}^\lambda$ 
13:    $\pi.blob, \pi.ttag \leftarrow \text{AE.ENC}(\pi.tk, \pi.seiv, \pi.blob)$ 
14:    $\psi \leftarrow \text{SE.SYMMKEYGEN}(\psi, \pi.h')$ 
15:   for  $\mathcal{K} \in \{iv, tag, tk, tiv, ttag\}$  do
16:      $\pi.\mathcal{K} \leftarrow \text{SE.SYMMENC}(\psi, \pi.h', \pi.tiv, \pi.\mathcal{K})$ 
17:   return  $\pi$ 
18:
19: procedure DEMS.DECRYPT( $\pi, \psi, \pi.pw$ )
20:   for  $\mathcal{K} \in \{iv, tag, tk, tiv, ttag\}$  do
21:      $\pi.\mathcal{K} \leftarrow \text{SE.SYMMDEC}(\psi, \pi.h', \pi.seiv, \pi.\mathcal{K})$ 
22:    $\pi.blob \leftarrow \text{AE.DEC}(\pi.tk, \pi.tiv, \pi.blob, \pi.ttag)$ 
23:   return  $\text{DESS.DECRYPT}(\pi, \psi, \pi.pw)$ 
24:
25: procedure DEMS.RATCHET( $\pi, \psi$ )
26:   for  $\mathcal{K} \in \{iv, tag, tk, tiv, ttag\}$  do
27:      $\pi.\mathcal{K} \leftarrow \text{SE.SYMMDEC}(\psi, \pi.h', \pi.seiv, \pi.\mathcal{K})$ 
28:    $\pi.blob \leftarrow \text{AE.DEC}(\pi.tk, \pi.tiv, \pi.blob, \pi.ttag)$ 
29:    $\pi.tk \leftarrow \text{AE.KEYGEN}(); \pi.tiv \xleftarrow{\$} \{0, 1\}^\lambda$ 
30:    $\pi.blob, \pi.ttag \leftarrow \text{AE.ENC}(\pi.tk, \pi.tiv, \pi.blob)$ 
31:    $\psi \leftarrow \text{SE.SYMMKEYGEN}(\psi, \pi.h'); \pi.seiv \xleftarrow{\$} \{0, 1\}^\lambda$ 
32:   for  $\mathcal{K} \in \{iv, tag, tk, tiv, ttag\}$  do
33:      $\pi.\mathcal{K} \leftarrow \text{SE.SYMMENC}(\psi, \pi.h', \pi.seiv, \pi.\mathcal{K})$ 
34:   return  $(\pi, \psi)$ 

```

---

**Security intuition.** This extra encryption layer is sufficient to withstand a multi-snapshot adversary as the adversary will always encounter a random-looking ciphertext. When the adversary finally gains access to the SE, they can only reverse the last outer encryption, but not any beforehand. A similar local technique of adding a reversible encryption layer is not possible without an SE as the encryption key would be part of the captured storage state. An alternative would be decrypting and re-encrypting the user data at every event. This would however be impractical as it necessitates the user to input their passphrase for each of these events.

### 5.3 Practical Implementation Details

Sloth is generally more practical than other schemes that rely on compute and memory-hard functions. This is because parameters of other schemes have to be chosen from the perspective of an attacker with access to a large parallel cluster of machines that are more powerful than a smartphone.

For the regular Sloth key stretching schemes and 1S-HiddenSloth, the implementation must not use a file system or underlying storage technology that allows an attacker to discover whether and when a file has been overridden. However, since the ratchet steps for MS-HiddenSloth are assumed to be predictable from the adversary's perspective, its security guarantees hold for any form of storage.

The actual implementation into a production app also needs to consider the allocated space for the deniable encryption since it should always be the same size regardless of the actual usage. Developers have to weigh the costs of having a large encrypted file for non-users against the storage requirements of its active users. Also, developers must take care that the deniable parts of the application do not leave any other traces in the form of crash reports and logging output.

In addition to application-specific storage, details on how the SE is implemented and stores state are important: limited on-chip storage means that SEs typically persist their state as encrypted key blobs and metadata using main device storage. While the techniques used by SEs generally provide strong confidentiality and integrity, they often only provide coarse roll-back protection. For instance, on Android, only OS updates cause the StrongBox roll-back protection counter to increment. As such, Android devices are vulnerable to roll-back attacks where the adversary can replace the encrypted version of the current state of SE storage<sup>4</sup> with an older version. The iOS documentation does not describe the level of roll-back protection provided on iPhones and iPads.

For our HiddenSloth scheme, which relies on ratcheting keys, the lack of roll-back protection can allow an adversary to detect changes if they both capture the encrypted key blobs and Sloth ciphertexts before and after a suspected usage event *and* later gain oracle access to the SE. However, we note that the encrypted key blobs are typically not accessible by application processes and they are never backed up by the system. Therefore, the ratcheting mechanism still provides value for common scenarios where an adversary gains access to stored application state via online backups that inadvertently include Sloth ciphertexts.

## 6 SECURITY ANALYSIS

We formally capture the security of our Sloth schemes using experiments between a challenger  $\mathcal{C}$  and a wall-time bounded (WT) adversary  $\mathcal{A}$  (Definition 2).

### 6.1 Key Stretching Security

We prove the security of LongSloth (Section 4.1) and RainbowSloth (Section 4.2). For this we formally define an SE-backed key stretching scheme and experiments for key stretching indistinguishability and hardness. The success of a WT adversary is controlled by the time required by the SE to execute operations  $T_{\text{SE-OP}}$  (and the Sloth parameters).

<sup>4</sup>Typically stored in `/data/misc/keystore/persistent.sqlite`



**Definition 6.** A key stretching scheme  $\Xi$  for fixed security parameter  $\lambda$  consists of two algorithms:

**KEYGEN** :  $\Psi \times \mathcal{P} \times H \rightarrow \Pi \times \Psi \times \{0, 1\}^\lambda$  takes an SE state, a password, and a key handle and returns a new storage state, an updated SE state, and the derived key.

**DERIVE** :  $\Pi \times \Psi \times \mathcal{P} \rightarrow \{0, 1\}^\lambda$  takes a storage state, SE state, and a password to derive a key  $k \in \{0, 1\}^\lambda$ .

We require for correctness that after the initialization operation  $(\pi, \psi, k) \leftarrow \text{KEYGEN}(\psi', pw, h)$  and all subsequent  $\text{DERIVE}(\pi, \psi, pw)$  executions return the same  $k$ .

We use the standard security notion for KDFs which requires “that the derived key [...] is indistinguishable from a random key if the adversary only knows the system parameter and the public [inputs]” [37]. Intuitively, this means that the output distribution meets typical requirements for input keys of other algorithms. We capture this using an experiment where the adversary must distinguish between a derived key and a random bit string.

**Definition 7** (Key Stretching Indistinguishability Experiment). Let  $\mathcal{P}$  be a probability distribution with Shannon entropy  $m$ . Let  $\mathcal{A}$  be a WT adversary and  $\mathcal{C}$  a WT challenger. Let  $\Xi$  be a key stretching scheme. Then the key stretching indistinguishability experiment  $\text{KEYIND}_{\mathcal{A}, \Xi}(\lambda, \mathcal{P})$  is defined as follows:

- (1)  $\mathcal{C}$  samples a password  $pw \xleftarrow{\$} \mathcal{P}$ . Let  $\psi$  be freshly initialized and  $h$  an arbitrary (but fixed) key handle.
- (2)  $\mathcal{C}$  computes  $\pi, \psi, k_0 \leftarrow \Xi.\text{KEYGEN}(\psi, pw, h)$  under  $\lambda$  and samples  $k_1 \leftarrow \{0, 1\}^\lambda$ .
- (3)  $\mathcal{C}$  randomly samples  $b \xleftarrow{\$} \{0, 1\}$ .
- (4)  $\mathcal{A}$  receives  $(\pi, \psi, k_b)$ .
- (5)  $\mathcal{A}$  receives oracle access  $O_{SE}$  (WT conditions).
- (6)  $\mathcal{A}$  outputs a bit  $b'$  and wins iff  $b = b'$ .
- (7) The experiment returns 1 iff  $\mathcal{A}$  wins, otherwise 0.

**Definition 8** (Key Stretching Indistinguishability). A key stretching scheme  $\Xi$  with operation cost  $\sigma$  is indistinguishable if for all WT adversaries  $\mathcal{A}$  with wall-time budget  $B$ , there is a negligible function  $\text{NEGL}$ :

$$\Pr[\text{KEYIND}_{\mathcal{A}, \Xi}(\lambda, \mathcal{P}) = 1] \leq \frac{1}{2} + \text{NEGL}(\lambda) + \frac{B}{\sigma \cdot 2^m},$$

where  $\mathcal{P}$  is a probability distribution with Shannon entropy  $m$ .

The LongSloth and RainbowSloth schemes respectively described in Section 4.1 and Section 4.2 fulfill these definitions as per the following theorems. We note that if  $m$  approaches  $\lambda$ , e.g. when the input is a strong cryptographic key, the definition can be updated so that the latter term is covered by  $\text{NEGL}(\lambda)$  which is similar to typical definitions. However, we are primarily interested in the settings where the password space is relatively small.

**Theorem 1** (LongSloth Indistinguishability). Let  $\mathcal{P}$  be a random distribution with Shannon entropy  $m$ . The key stretching scheme LongSloth is indistinguishable.

**PROOF.** We present a full proof in Appendix C.1.2.  $\square$

**Theorem 2** (RainbowSloth Indistinguishability). Let  $\mathcal{P}$  be a random distribution with Shannon entropy  $m$ . The key stretching scheme RainbowSloth is indistinguishable.

**PROOF.** We present a full proof in Appendix C.2.2.  $\square$

In addition to the theoretical result above, we are interested in the practical success rate of a brute-force attacker, i.e., given  $k$  find  $pw$ , with wall time budget  $B$ .

**Definition 9** (Key Stretching Hardness Experiment). Let  $\mathcal{P}$  be a probability distribution with Shannon entropy  $m$ . Let  $\mathcal{A}$  be a WT adversary and  $\mathcal{C}$  a WT challenger. Let  $\Xi$  be a key stretching scheme. Then the key stretching hardness experiment  $\text{KEYHARD}_{\mathcal{A}, \Xi}(\lambda, \mathcal{P})$  is defined below:

- (1)  $\mathcal{C}$  samples a password  $pw \xleftarrow{\$} \mathcal{P}$ . Let  $\psi$  be freshly initialized and  $h$  an arbitrary (but fixed) key handle.
- (2)  $\mathcal{C}$  computes  $k = \Xi.\text{KEYGEN}(\psi, pw, h)$ .
- (3)  $\mathcal{A}$  receives  $(\pi, \psi, k)$ .
- (4)  $\mathcal{A}$  receives oracle access  $O_{SE}$  (WT conditions).
- (5)  $\mathcal{A}$  outputs a  $pw'$  and wins iff  $k = \Xi.\text{DERIVE}(\psi, pw', h)$ .
- (6) The experiment returns 1 iff  $\mathcal{A}$  wins, otherwise 0.

**Definition 10** (Key Stretching Hardness). A key stretching scheme  $\Xi$  is  $\sigma$ -hard if for all WT adversaries  $\mathcal{A}$  with wall-time budget  $B$ ,

$$\Pr[\text{KEYHARD}_{\mathcal{A}, \Xi}(\lambda, \mathcal{P}) = 1] \leq \frac{B}{\sigma \cdot 2^m},$$

where  $\mathcal{P}$  is a probability distribution with Shannon entropy  $m$ .

**Theorem 3** (LongSloth Hardness). Let  $\mathcal{P}$  be a random distribution with Shannon entropy  $m$ . Then the key stretching scheme LongSloth with parameter  $l$  is  $(l \cdot c_{\text{HMAC}})$ -hard.

**PROOF.** We present a full proof in Appendix C.1.1.  $\square$

**Theorem 4** (RainbowSloth Hardness). Let  $\mathcal{P}$  be a random distribution with Shannon entropy  $m$ . Then the key stretching scheme RainbowSloth with parameter  $n$  is  $(n \cdot c_{\text{ECDH}})$ -hard.

**PROOF.** We present a full proof in Appendix C.2.1.  $\square$

## 6.2 Deniable Encryption Security

Similarly to key stretching, we discuss the security of deniable encryption by giving formal definition, describing security experiments, and then showing the security of HiddenSloth. We only discuss the multi-snapshot case and prove the security of MS-HiddenSloth. A similar (and simpler) reasoning applies to 1S-HiddenSloth.

**Definition 11** (SE with MS Deniable Encryption Support). A SE-backed deniable encryption scheme  $\Delta$  for fixed security parameter  $\lambda$  and maximum data size  $s \in \mathbb{N}$  consists of three algorithms:

**INIT** :  $\Psi \times H \times \mathbb{N} \rightarrow \Pi \times \Psi$  takes a SE state, a key handle, and storage size  $s$  and returns a new storage state and SE state.

**ENCRYPT** :  $\Pi \times \Psi \times \mathcal{P} \times \{0, 1\}^s \rightarrow \Pi$  updates a storage state for given SE state and a password so that it now stores the given data.

**DECRYPT** :  $\Pi \times \Psi \times \mathcal{P} \rightarrow \{0, 1\}^s \cup \{\perp\}$  tries to decrypt from the storage state given the password. If successful, the previously stored data is returned, otherwise  $\perp$ .

**RATCHET** :  $\Pi \times \Psi \rightarrow \Pi \times \Psi$  updates the storage and SE state by decrypting and re-encrypting the stored data.

In our indistinguishability experiment the adversary interactively builds any two valid sequences of starting with an `INIT` operation and multiple `ENCRYPT` operations. After each encryption, they receive the persisted states. The challenger will then call `RATCHET` on one of the sequences and provide the new state. The inability of the adversary to decide which of the sequences the new state has been derived from, implies that the `RATCHET` operation is indistinguishable from an `ENCRYPT` operation—hence the scheme provides effective security against a multi-snapshot adversary.

**Definition 12** (MS Deniable Encryption Indistinguishability Experiment). Let  $\lambda$  be a fixed security parameter and  $s$  the maximum data size. Let  $\mathcal{P}$  be a probability distribution with Shannon entropy  $m$ . Let  $\mathcal{A}$  be a WT adversary with wall time budget  $B$  and  $C$  a WT challenger. Let  $\Delta$  be a multi-snapshot deniable encryption scheme. Then the multi-snapshot deniable encryption indistinguishability experiment  $\text{DE-MS-IND}_{\mathcal{A},\Delta}(\lambda, \mathcal{P})$  is defined as follows:

- (1)  $C$  samples  $pw \xleftarrow{\$} \mathcal{P}$ . Let  $\psi_0, \psi_1$  be freshly initialized and  $h$  an arbitrary, but fixed key handle.
- (2)  $C$  randomly samples  $b \xleftarrow{\$} \{0, 1\}$
- (3)  $C$  computes  $\pi_0, \psi_0 \leftarrow \Delta.\text{INIT}(\psi_0, h, s)$  and  $\pi_1, \psi_1 \leftarrow \Delta.\text{INIT}(\psi_1, h, s)$  under  $\lambda$ .
- (4)  $A$  sends any  $\tilde{b} \in \{0, 1\}$  and  $m \in \{0, 1\}^s$  to  $C$ .
- (5)  $C$  executes  $\pi_{\tilde{b}} \leftarrow \Delta.\text{ENCRYPT}(\pi_{\tilde{b}}, \psi_{\tilde{b}}, pw, m)$ , then  $C$  sends  $\pi_{\tilde{b}}$  to  $A$ .
- (6)  $A$  can repeat execution from step 4 several times (under WT conditions).
- (7)  $C$  randomly samples  $b' \xleftarrow{\$} \{0, 1\}$ , computes  $\pi_b, \psi_b \leftarrow \Delta.\text{RATCHET}(\pi_b, \psi_b)$ , and provides  $A$  with  $(\pi_b, \psi_b)$ .
- (8)  $A$  receives oracle access  $O_{SE}$  under the WT conditions.
- (9)  $A$  outputs a bit  $b'$  and wins iff  $b = b'$ .
- (10) The experiment returns 1 iff  $A$  wins, otherwise 0.

**Definition 13** (MS Deniable Encryption Indistinguishability). A multi-snapshot deniable encryption scheme  $\Delta$  instantiated with a  $\sigma$ -hard key stretching scheme is  $m$ -entropy secure if for all WT adversaries  $\mathcal{A}$  with time budget  $B$ , there is a function  $\text{NEGL}$  such that

$$\Pr[\text{DE-MS-IND}_{\mathcal{A},\Delta}(\lambda, \mathcal{P}) = 1] \leq \frac{1}{2} + \text{NEGL}(\lambda) + \frac{B}{\sigma \cdot 2^m},$$

where  $\mathcal{P}$  is a probability distribution with Shannon entropy  $m$ .

**Theorem 5.** (MS-HiddenSloth Indistinguishability) The multi-snapshot deniable encryption scheme MS-HiddenSloth instantiated with a  $\sigma$ -hard key stretching scheme is  $m$ -entropy secure.

**PROOF.** We present a full proof in Appendix C.3.  $\square$

In particular, this means that MS-HiddenSloth instantiated with our key stretching schemes LongSloth or RainbowSloth is  $m$ -entropy secure. Similarly to the key stretching scheme, we also analyze the practical success rate of a brute-force attacker.

**Definition 14** (Deniable Encryption Hardness Experiment). Let  $\mathcal{P}$  be a probability distribution with Shannon entropy  $m$ . Let  $\mathcal{A}$  be a WT adversary with wall time budget  $B$  and  $C$  a WT challenger. Let  $\Delta$  be a key stretching scheme. Then the deniable encryption hardness experiment  $\text{DEHARD}_{\mathcal{A},\Delta}(\lambda, \mathcal{P})$  is defined as follows:

- (1)  $\mathcal{A}$  provides  $C$  with  $data$  with  $|data| > 0$ .
- (2)  $C$  samples a password  $pw \xleftarrow{\$} \mathcal{P}$ . Let  $\psi$  be freshly initialized and  $h$  an arbitrary (but fixed) key handle. Let  $\pi, \psi = \Delta.\text{INIT}(\psi, h, |data|)$ .
- (3)  $C$  computes  $\pi = \Delta.\text{ENCRYPT}(\pi, \psi, pw, data)$ .
- (4)  $\mathcal{A}$  receives  $(\pi, \psi)$ .
- (5)  $\mathcal{A}$  receives oracle access  $O_{SE}$  (WT conditions).
- (6)  $\mathcal{A}$  outputs a  $pw'$  and wins iff  $data = \Delta.\text{DECRYPT}(\pi, \psi, pw')$ .
- (7) The experiment returns 1 iff  $\mathcal{A}$  wins, otherwise 0.

**Definition 15** (Deniable Encryption Hardness). A deniable encryption scheme  $\Delta$  is  $\sigma$ -hard if for all WT adversaries  $\mathcal{A}$  with wall-time budget  $B$ ,

$$\Pr[\text{DEHARD}_{\mathcal{A},\Delta}(\lambda, \mathcal{P}) = 1] \leq \frac{B}{\sigma \cdot 2^m},$$

where  $\mathcal{P}$  is a probability distribution with Shannon entropy  $m$ .

**Theorem 6** (MS-HiddenSloth Hardness). Let  $\mathcal{P}$  be a random distribution with Shannon entropy  $m$ . Then the deniable encryption scheme MS-HiddenSloth instantiated with a  $\sigma$ -hard key stretching scheme is  $\sigma$ -hard.

**PROOF.** We present a full proof in Appendix C.3.2.  $\square$

## 7 EVALUATION

We measure the performance of SEs in Android and iOS devices (Section 7.1), choose practical parameters for our schemes (Section 7.2), and test full implementations of LongSloth and RainbowSloth (Section 7.3) as well as HiddenSloth (Section 7.4).

### 7.1 Performance Characteristics of SEs

For iOS devices we measure the duration of the Secure Enclave's ECDH operations on the recent iPhones from XR to 14. These cover the Apple chips A12, A13, A14, and A15. We wrote a benchmark app that first creates a secret P-256 key within the SE. It then creates a new random public key and performs ECDH with the private key within the SE. We repeat the ECDH step 1,000 times with the new public keys and measure the elapsed time. The results (Figure 3) show that similar chips have similar run times, e.g., A15 for iPhone 13 and iPhone 14. All measurements are between 6 ms and 16 ms with little variance per model. Surprisingly, the A13 chip has lower throughput than its predecessor. However, A13 is also the first one with a mathematically verified public key implementation [4] which might point to a difference in the underlying algorithm.

For Android devices, we measure the duration of HMAC executions for inputs of varying lengths on multiple devices. We wrote a benchmark tool that first creates a secret key within the SE. It then generates random byte arrays as inputs, passes these as input to the SE, and receives the computed HMAC value as a result. We repeat this step 10 times for each device and input length to measure the elapsed time. The results are shown in Figure 4. Intuitively, the measured time increases with input length in an almost linear manner. There are minor inflection points that differ between the devices. In our experiments, the Samsung phones also have a 10× higher bandwidth compared to the Google devices. For our parameter choice, this will result in a larger  $l$  value for those devices. Notably, for the

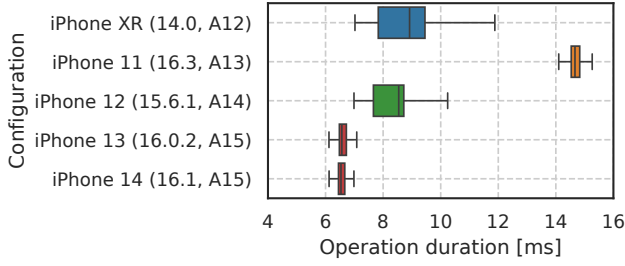


Figure 3: The duration of the ECDH operation on iOS for different phones (iOS version and chip generation in brackets).

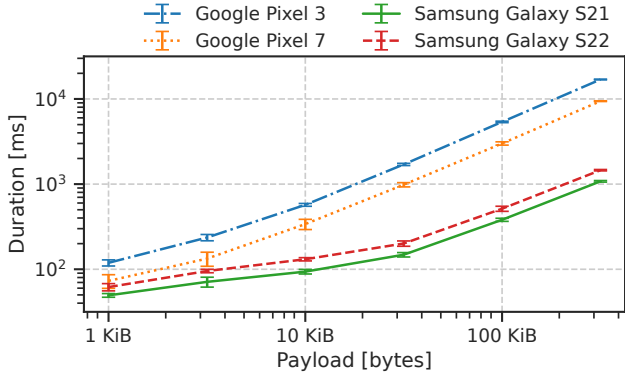


Figure 4: The duration of a HMAC operations on Android phones with StrongBox support. Both axes are log-scale.

	Entropy	50 years	100 years
WordList (3 words)	38.8	3.4 ms	6.7 ms
WordList (4 words)	51.7	< 0.1 ms	< 0.1 ms
AlphaNum (5 chars)	29.8	1.7 s	3.4 s
AlphaNum (6 chars)	35.7	27.8 ms	55.5 ms
AlphaNum (7 chars)	41.7	0.4 ms	0.9 ms
PIN (6 digits)	19.9	26.3 min	52.6 min
PIN (9 digits)	29.9	1.6 s	3.2 s

Table 2: Overview of password configurations, their entropy (bits), and the required  $t_{individual}$  for given security margins.

same device and input length, the variance is very small. Once the parameters have been established for each device, their impact is predictable and dependable.

## 7.2 Parameter choice

For our evaluation, we target a security level of  $T_{total} = 100$  years. We treat the acceptable password complexity, such as the alphabet size and the number of characters, as input parameters. Other than

classic alphanumerical<sup>5</sup> passwords, we also consider passphrases that are based on the EFF word list [16] that contains 7776 words and PINs consisting of only digits. Let  $|A|$  be the alphabet size and  $|pw|$  the password length, then the entropy for a configuration is  $e = \log_2(|A|^{|pw|})$ . We want that an attacker’s worst-case to match the targeted security level, i.e., if they brute-force the entire space, then they require at least  $T_{total}$ . Let  $T_{total}$  be the targeted security level, then for a password configuration with entropy  $e$  each password verification should take at least  $t_{individual} = T_{total} \cdot 2^{-e}$ . Table 2 shows that while short passphrases and alpha-numerical passwords are feasible, short PIN codes are not as they require impractically long  $t_{individual}$  times.

We think that using  $T_{total}$ , i.e. the time to search the entire passphrase space, is convenient for interpretation of our results. This is because its linearity allows to calculate the expected time for a given adversary success rate (and the other way around). For example, for our choice of  $T_{total} = 100$  years, we would expect that after a 10 years the adversary has guessed the correct passphrase with a 10% chance. If we would like to reduce that chance to 1%, we must set  $T_{total} = 1000$  years and update  $t_{individual}$  accordingly. Particular attention must be paid in scenarios where the adversary succeeds when guessing 1 out of  $n$  passphrases, e.g. after capturing smartphones from a group of people. Since the individual phones can be brute-forced in parallel, the random variables are independent and the overall success rate increases accordingly.

The password configurations that we show in Table 2 are shorter than those used in web applications nowadays. This is because for Sloth the parameter choice does not need to conservatively assume an attacker with highly parallel computing resources. Short passwords, and in particular passphrases generated from word lists, are more memorable and can be generated by the application for the user – and thus avoid the problem of users picking weak passwords or reusing the same one for different services.

In our main evaluation (Section 7.3) we will examine two configurations: 3-word passphrases from the EFF word list ( $c1$ ) and 6-character alphanumerical passwords ( $c2$ ). In both cases, we multiply the minimal  $t_{individual}$  (as per Table 2) by a **safety factor of  $\times 10$**  for conservative parameter choice. Therefore:  $t_{c1} = 67$  ms and  $t_{c2} = 555$  ms. This safety factor can account for the attacker over-clocking the SE and overhead by the operating system when communicating with the chip that an attacker might be able to “optimize away”.

With the required minimum times  $t_{c1}, t_{c2}$  for the individual operations, we determine  $l$  for LongSloth on Android and  $n$  for RainbowSloth on iOS. On Android, we use the measurements from Figure 4 to fit a second-degree polynomial where we set the y-values to the smallest measurement for a given size minus 2 standard deviations. We then pick the  $l$  value at the intersection with the desired duration and round to the next multiple of 100. On iOS we pick the 10th percentile value  $t_{p10}$  of our measurements as a conservative worst-case (i.e. fastest) duration and then compute  $n_c = \lceil \frac{t_c}{t_{p10}} \rceil$  for  $c \in \{c1, c2\}$ . The resulting values are summarized in Table 3. If an app cannot find existing parameters for a new device type, it can perform a similar method on its first start to self-calibrate.

<sup>5</sup>Case-sensitive letters a–zA–Z and digits 0–9, hence  $|A| = 62$ .

	3 words $t_{c1} = 67 \text{ ms}$	6 characters $t_{c2} = 555 \text{ ms}$
LongSloth (parameter: l)		
Google Pixel 3	1,500	11,600
Google Pixel 7	4,500	24,200
Samsung Galaxy S21	10,700	178,000
Samsung Galaxy S22	2,200	145,100
RainbowSloth (parameter: n)		
iPhone 11	6	43
iPhone 12	10	76
iPhone 13	12	95

**Table 3: Parameter choices for given password configuration and its  $t_{individual}$  augmented by a safety factor of  $10\times$ .**

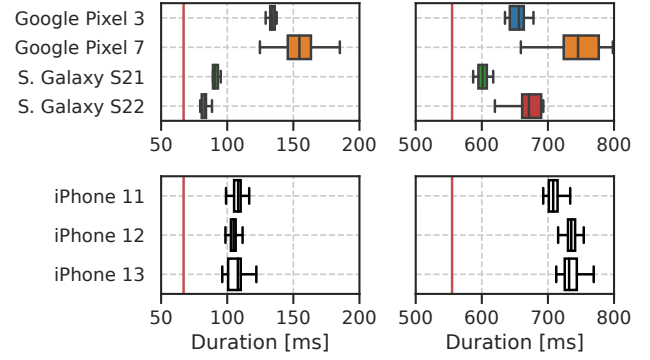
### 7.3 LongSloth and RainbowSloth

We implemented LongSloth on Android and RainbowSloth on iOS. All code including documentation and analysis scripts is available under an MIT license.<sup>6</sup> Where possible we use the existing cryptography APIs of the platform, with AES-GCM for symmetric authenticated encryption, and HKDF-SHA256 as the KDF. For PwHASH we use the third-party LibSodium library which is available on both platforms and implements the memory-hard password hashing algorithm Argon2id [14]. For the evaluation, we choose the recommended OWASP parameters for Argon2id with 19 MiB of memory and an iteration count of 2 [32]. Since one could opt for another PwHASH implementation, we exclude its runtime (50 ms) from our results for LongSloth and RainbowSloth.

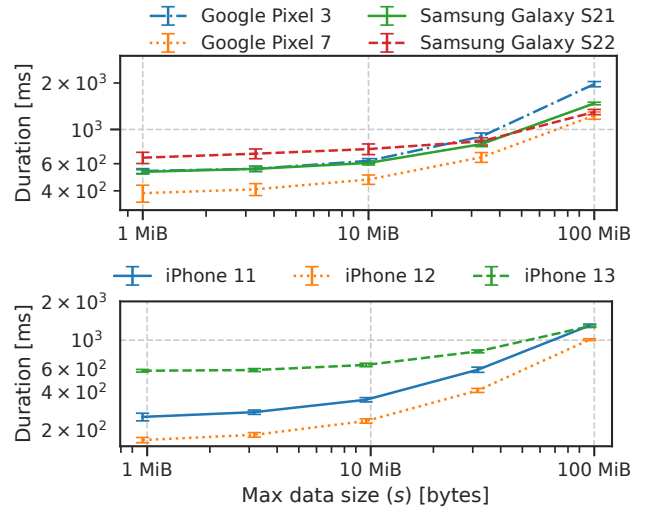
For our evaluation, we are interested in the duration of the SLOTH.DERIVE operations, as these determine the time costs for an attacker. We use the parameters chosen in Section 7.2 including the safety factor and execute each operation 10 times. The results are shown in Figure 5 for LONGSLOTH on Android RAINBOWSLOTH on iOS. In all cases, the measured total durations comfortably exceeded the threshold times  $t_{c1}$  and  $t_{c2}$ . This confirms that with our parameter choice, the algorithm meets its minimum timing promise and hence key stretching security. The variance for the individual configurations is small which can allow for reducing the safety factor.

### 7.4 HiddenSloth

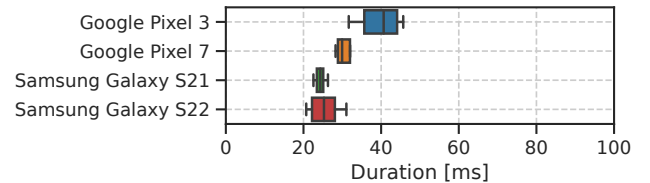
For the deniable encryption scheme HiddenSloth we evaluate its RATCHET methods for varying maximum data sizes  $s$ . This parameter  $s$  spans from a small storage size that might wrap text-only configurations (1 MiB) to larger ones that can store long chat histories including media (100 MiB). Our results for Android and iOS are shown in Figure 6. The measured durations are similar for all test devices and range from around 700 ms for 1 MiB to about 2 s for storage with 100 MiB capacity. We note that the RATCHET step does not require the password and thus can be executed in the



**Figure 5: The duration of the LONGSLOTH.DERIVE operation on Android (top) RAINBOWSLOTH.DERIVE operation on iOS (bottom). For both we evaluated the configurations  $c1$  (left) and  $c2$  (right) from Table 3 for different phones. The red lines indicate the threshold times  $t_{c1}$  and  $t_{c2}$ .**



**Figure 6: Duration of the HIDDEN SLOTH.RATCHET step for various max size  $s$  on Android (top) and iOS (bottom). Both axes are log-scale.**



**Figure 7: Duration of HIDDEN SLOTH random-access decryption without authentication for 1 MiB blocks.**

<sup>6</sup><https://github.com/lambdaioneer/sloth>

background without any user interaction. As such, even large storage sizes have no user-visible impact and we suggest scheduling it when the device is idle and charging.

Decryption speed is critical for most application use-cases as it dictates the speed by which stored data, e.g. photos, load. Apps can optimize this step by caching the derived key  $k$  and unwrapped  $\pi.K$  in memory after the user has entered their password. For this we implemented authenticated encryption using an encrypt-then-mac regime with AES-CTR and HMAC. This then allows random access to individual blocks of the ciphertext. We evaluate random-access decryption on Android and the results in Figure 7 show that reading a 1 MiB block generally takes less than 50 ms. The app should verify the authentication tag of the entire storage before performing decryption operations.

## 7.5 Limitations

Sloth remains vulnerable to physical attacks during usage such as shoulder surfing [10], smudge attacks [3], and side-channel attacks that monitor keyboard entry [40]. However, this is compatible with our threat model and solutions for these are orthogonal to our scheme. In addition, if an app developer allows user-chosen passphrases, then our Shannon entropy assumption (Section 3.1) is not valid. Therefore, we designed Sloth such that randomly-chosen passphrases are feasible. Since we rely on a hardware-backed secret, device loss means that the secret cannot be recovered. However, apps can introduce complementary backup procedures<sup>7</sup>.

## 8 RELATED WORK

The terminology of key stretching which we also use for our paper was coined in 1999 by Kelsey et al. [27]. Examples of modern password hashing functions include PBKDF2 [24], scrypt [33], and Argon2 [14]. All share the property that an attacker can increase brute-force speed by using more computers. Blocki et al. [15] show in an economic analysis of offline password cracking that key stretching alone does not provide sufficient protection.

Password-authenticated key agreement (PAKE) protocols, such as OPAQUE [23], use interactive protocols with a server to rate-limit password guesses. However, such protocols are not suitable when there is unreliable Internet connectivity. Further, the required network communication leaves traces that rule out their usage for deniable encryption schemes. This also holds for related protocols such as DupLESS [26]. The traffic can be obfuscated by sending “dummy” requests. However, on smartphone mobile network transmission have high energy costs and thus limiting how often this can be done.

The VeraCrypt project [36] is the most-popular encryption software that offers support plausible deniable encryption for regular computers. It can also boot into different operating systems depending on the entered password when unlocking the machine. To our knowledge, StegFS [30] was the first practical design of a plausibly-deniable encrypted file system. The recent INVISILINE design [34] encodes a hidden volume in the IVs of a general-purpose

block device encryption scheme and provides multi-snapshot resistance. Our Sloth key stretching algorithm can be used as a drop-in replacement with any of these schemes to improve their resistance to brute-force attacks.

For mobile devices, MobiFlage [41] describes a design for providing deniable encryption on Android resisting single-snapshot adversaries. The MobiCeal [18] paper describes an implementation that protects against multi-snapshot adversaries by obfuscating access patterns using “dummy writes”. However, both MobiFlage and MobiCeal require changes to the operating system which renders them impractical for use by developers who wish to support users of standard handsets. The work by Liao et al. [29] describes a theoretical design where some computations for the deniable encryption scheme are run inside an ARM TrustZone environment. Likewise, this design would require significant changes to the operating system and the ability to run custom code within the TEE. Also, none of the above protect against brute-force attacks.

Our work is inspired by a deniable encryption scheme sketched in CoverDrop [1] and shares the idea of storing a separate key inside an SE. However, CoverDrop does not allow for generic key stretching and it requires the entire ciphertext to grow to guarantee lower guessing rates. Therefore, the security parameters cannot be independently tweaked and it is less efficient than Sloth. Furthermore, the approach described in CoverDrop does not work on iOS due to the lack of AES-GCM support in Apple’s SE. The paper also does not offer security proofs nor practical evaluations for a representative set of devices.

Eldridge et al. [19] implement support for one-time programs on smartphones. Similar to us, they creatively work around the limitations of embedded SEs. Their solution allows creating one-time programs for general functionalities using only a counter box primitive. However, such counter boxes are not exposed to end-user apps on iOS and thus they require jailbroken devices.

## 9 CONCLUSION

Sloth is a family of novel practical key stretching and deniable encryption schemes which leverages the limited throughput of the SE to provide strong security guarantees and favorable parameter choices for standard smartphones. Sloth works without any changes to the operating system and allows for shorter passphrases as the adversary cannot speed up brute-force attacks by using multiple computers. Our survey shows that SEs are more widely available and more practical than is generally assumed. We presented a formal model of SEs which capture their timing characteristics. Instead of asymptotic analysis, we work with absolute duration to provide definite time bounds. This precision allows picking smaller parameters that improve efficiency and usability for the user. We believe that SEs will play a critical role in the security of mobile applications and services, and hope to inspire more research in the area of hardware-assisted security on smartphones.

## ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers from this and previous submissions for their valuable feedback. Daniel Hugenroth is supported by Nokia Bell Labs. This work is partially funded by Mysten Labs.

<sup>7</sup>For example, by adding an additional layer of indirection where the derived secret  $k$  encrypts another secret  $k'$  that is then used, e.g., for HIDDEN SLOTH. The  $k'$  can be temporarily decrypted, transferred to another device, and then re-encrypted there using a new derived secret  $k_2$ .

## REFERENCES

- [1] Mansoor Ahmed-Rengers, Diana A Vasile, Daniel Hugenroth, Alastair R Beresford, and Ross Anderson. 2022. CoverDrop: Blowing the Whistle Through A News App. *Proceedings on Privacy Enhancing Technologies* 2022, 2 (2022), 47–67.
- [2] Amazon. 2023. AWS Device Farm. <https://aws.amazon.com/device-farm>.
- [3] MD Amruth and K Praveen. 2016. Android smudge attack prevention techniques. In *Intelligent Systems Technologies and Applications: Volume 2*. Springer, 23–31.
- [4] Apple Inc. 2021. Apple Platform Security - Secure Enclave. <https://support.apple.com/en-gb/guide/security/sec59b0b31ff/web>.
- [5] Apple Inc. 2022. App Store - iOS and iPadOS usage. <https://web.archive.org/web/20230130101543/https://developer.apple.com/support/app-store/>.
- [6] Apple Inc. 2022. Apple Security Bounty. <https://security.apple.com/bounty/>.
- [7] Apple Inc. 2023. Apple Developer Documentation - SecureEnclave. <https://developer.apple.com/documentation/cryptokit/secureenclave>.
- [8] Apple Inc. 2023. Apple Developer Documentation - SecureEnclave.P256. <https://developer.apple.com/documentation/cryptokit/secureenclave/p256>.
- [9] Diego F Aranha, Pierre-Alain Fouque, Chen Qian, Mehdi Tibouchi, and Jean-Christophe Zapolowicz. 2014. Binary elligator squared. In *International Conference on Selected Areas in Cryptography*. Springer, 20–37.
- [10] Adam J Aviv, John T Davin, Flynn Wolf, and Ravi Kuber. 2017. Towards baselines for shoulder surfing on mobile authentication. In *33rd Annual Computer Security Applications Conference*. 486–498.
- [11] Feng Bao, Robert H Deng, and Huafei Zhu. 2003. Variations of diffie-hellman problem. In *Information and Communications Security: 5th International Conference, ICICS 2003, Huhehaote, China, October 10-13, 2003. Proceedings 5*. Springer, 301–312.
- [12] Mihir Bellare. 2006. New proofs for NMAC and HMAC: Security without collision-resistance. In *Crypto*, Vol. 4117. Springer, 602–619.
- [13] Daniel J Bernstein, Mike Hamburg, Anna Krasnova, and Tanja Lange. 2013. Elligator: elliptic-curve points indistinguishable from uniform random strings. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*. 967–980.
- [14] Alex Biryukov, Daniel Dinu, and Dmitry Khovratovich. 2016. Argon2: new generation of memory-hard functions for password hashing and other applications. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 292–302.
- [15] Jeremiah Blocki, Benjamin Harsha, and Samson Zhou. 2018. On the economics of offline password cracking. In *2018 IEEE Symposium on Security and Privacy (S&P)*. IEEE, 853–871.
- [16] Joseph Bonneau. 2016. Deep Dive: EFF’s New Wordlists for Random Passphrases. *Electronic Frontier Foundation (EFF)* (2016). <https://www.eff.org/deeplinks/2016/07/new-wordlists-random-passphrases> Accessed September 2023.
- [17] Daniel R. L. Brown. 2009. *SEC 1: Elliptic Curve Cryptography*. Standard. Certicom Research. Version 2.0.
- [18] Bing Chang, Fengwei Zhang, Bo Chen, Yingjiu Li, Wen-Tao Zhu, Yangguang Tian, Zhan Wang, and Albert Ching. 2018. Mobiceal: Towards secure and practical plausibly deniable encryption on mobile devices. In *2018 48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 454–465.
- [19] Harry Eldridge, Aarushi Goel, Matthew Green, Abhishek Jain, and Maximilian Zinkus. 2022. One-Time Programs from Commodity Hardware. In *Theory of Cryptography Conference*. 121–150.
- [20] Google Inc. 2022. Hardware Security Best Practices. <https://source.android.com/docs/security/best-practices/hardware>.
- [21] Google Inc. 2022. Android and Google Devices Security Reward Program Rules. <https://bughunters.google.com/about/rules/6171833274204160/android-and-google-devices-security-reward-program-rules>.
- [22] Google Inc. 2022. Android Keystore system - Hardware security module. <https://developer.android.com/training/articles/keystore#HardwareSecurityModule>.
- [23] Stanislaw Jarecki, Hugo Krawczyk, and Jiayu Xu. 2018. OPAQUE: an asymmetric PAKE protocol secure against pre-computation attacks. In *Advances in Cryptology—EUROCRYPT*. Springer, 456–486.
- [24] Burt Kaliski and A Rusch. 2017. RFC 8018: PKCS # 5: Password-Based Cryptography Specification Version 2.1.
- [25] Jonathan Katz and Yehuda Lindell. 2020. *Introduction to modern cryptography*. CRC press.
- [26] Sriram Keelveedhi, Mihir Bellare, and Thomas Ristenpart. 2013. DupLESS: Server-Aided encryption for deduplicated storage. In *22nd USENIX Security Symposium (USENIX Security '13)*. 179–194.
- [27] John Kelsey, Bruce Schneier, Chris Hall, and David Wagner. 2005. Secure applications of low-entropy keys. In *Information Security: First International Workshop, ISW'97 Tatsunokuchi, Ishikawa, Japan September 17–19, 1997 Proceedings*. Springer, 121–134.
- [28] Hugo Krawczyk. 2010. Cryptographic Extraction and Key Derivation: The HKDF Scheme.. In *CRYPTO*, Vol. 6223. Springer, 631–648.
- [29] Jinghui Liao, Bo Chen, and Weisong Shi. 2021. TrustZone enhanced plausibly deniable encryption system for mobile devices. In *2021 IEEE/ACM Symposium on Edge Computing (SEC)*. IEEE, 441–447.
- [30] Andrew D McDonald and Markus G Kuhn. 2000. StegFS: A steganographic file system for Linux. In *Information Hiding: Third International Workshop, IH'99, Dresden, Germany, September 29-October 1, 1999 Proceedings 3*. Springer, 463–477.
- [31] Alexander Nilsson, Pegah Nikbakht Bideh, and Joakim Brorsson. 2020. A survey of published attacks on Intel SGX. *arXiv preprint arXiv:2006.13598* (2020).
- [32] Open Worldwide Application Security Project (OWASP). 2021. Password storage cheat sheet. [https://cheatsheetseries.owasp.org/cheatsheets/Password\\_Storage\\_Cheat\\_Sheet.html](https://cheatsheetseries.owasp.org/cheatsheets/Password_Storage_Cheat_Sheet.html) Accessed September 2023.
- [33] Colin Percival and Simon Josefsson. 2016. The scrypt password-based key derivation function. <https://www.rfc-editor.org/rfc/rfc7914.html>.
- [34] Sandeep Kiran Pinjala, Bogdan Carbutar, Anrini Chakraborti, and Radu Sion. 2023. INVISILINE: Invisible Plausibly-Deniable Storage. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 18–18.
- [35] Sandro Pinto and Nuno Santos. 2019. Demystifying Arm TrustZone: A comprehensive survey. *ACM computing surveys (CSUR)* 51, 6 (2019), 1–36.
- [36] VeraCrypt project. 2023. VeraCrypt - Free Open source disk encryption with strong security for the Paranoid. <https://www.veracrypt.fr/en/Home.html>.
- [37] Baodong Qin, Shengli Liu, Tsz Hon Yuen, Robert H Deng, and Kefei Chen. 2015. Continuous non-malleable key derivation and its application to related-key security. In *IACR International Workshop on Public Key Cryptography*. Springer, 557–578.
- [38] Keegan Ryan. 2019. Hardware-backed heist: Extracting ECDSA keys from Qualcomm’s TrustZone. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 181–194.
- [39] Alon Shalevsky, Eyal Ronen, and Avishai Wool. 2022. Trust Dies in Darkness: Shedding Light on Samsung’s TrustZone Keymaster Design. *IACR Cryptol. ePrint Arch.* 2022 (2022), 208.
- [40] Laurent Simon, Wenduan Xu, and Ross Anderson. 2016. Don’t Interrupt Me While I Type: Inferring Text Entered Through Gesture Typing on Android Keyboards. In *Proceedings on Privacy Enhancing Technologies*, Vol. 3. 136–154.
- [41] Adam Skillen and Mohammad Mannan. 2013. Mobiflage: Deniable storage encryption for mobile devices. *IEEE Transactions on Dependable and Secure Computing* 11, 3 (2013), 224–237.
- [42] Mehdi Tibouchi. 2014. Elligator squared: Uniform points on elliptic curves of prime order as uniform random strings. In *International Conference on Financial Cryptography and Data Security*. Springer, 139–156.

## A SURVEY OF SE AVAILABILITY

The availability of SE functionality for end-user apps is determined by hardware support (i.e. does the device have an SE) as well as platform support for the respective API. We estimate market share for both Apple’s iOS devices and Google’s Android ecosystem.

## A.1 SE Support on Apple devices

As discussed in Section 2, Apple first added Secure Enclaves to their iPhone 5S in 2013 and an API was added in iOS 13 which was released in 2019. Since the iPhone 5S only supported iOS version until 12, all devices with iOS 13 (platform support) also contain a Secure Enclave (hardware support). This is supported by the fact that there is no API to check for the presence of a Secure Enclave. The App Store statistics for May 2022 show that 82% of all iPhones use iOS 15 and 14% use iOS 14 [5]. No data is given for iOS 13 or before. Therefore, at least 96% of all iPhones devices expose SE functionality to developers and are compatible with our Sloth scheme.

## A.2 SE Support on Android devices

The situation for the Android ecosystem is more complex as devices are manufactured by different vendors with different hardware as well as changes to the operating system. This can lead to situations where Android devices have an SE, but do not offer API access, or where a device’s API is compatible, but has no SE. We use the overall distribution of active Android versions as an upper boundary for platform support. We then execute test code on dozens of real devices to determine hardware support.



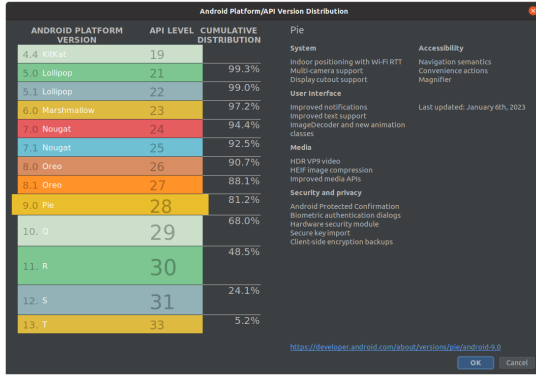


Figure 8: Distribution of Android version as shown in the Android Studio IDE for January 6th, 2023.

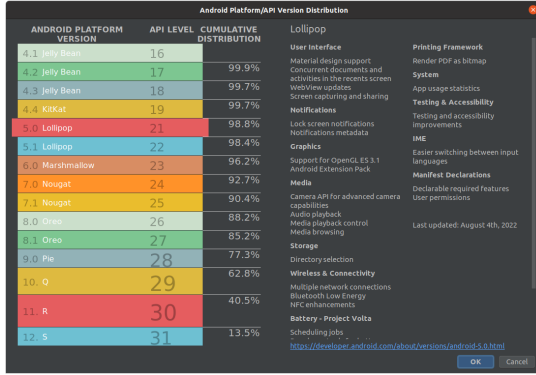


Figure 9: Distribution of Android version as shown in the Android Studio IDE for August 4th, 2022.

For the platform support we use the API Version Distribution that is shown in the “New Project” wizard in the Android Studio IDE. Figure 8 shows the API distribution that is shown when creating new project in the Android Studio IDE when it was last updated on January 6th, 2023. Figure 9 shows the same data for the previous update in August 4th, 2022. Figure 8 shows that 97.2% of devices run API level 23 (Android M) or higher and hence provide the API to check whether a key is backed by a TEE or SE. Furthermore, 81.2% of devices run API level 28 (Android P) or higher and hence provide the StrongBox API that can enforce storage in an SE.

For hardware support, we use the AWS Device Farm [2] which allows remote access to real devices in a datacenter. We compile a test application and upload it for execution on all selected device types. Our test code performs multiple checks. First, we create different key types and then check via `KeyInfo#isInsideSecureHardware` if it is stored in secure hardware. A positive answer indicates that the device has a TEE or SE. Then, we read the `PackageManager` feature flag for StrongBox support: `FEATURE_STRONGBOX_KEYSTORE`. Finally, we create different key types with the `setIsStrongBoxBacked` property that enforces storing them in an SE and checking for failures. Figure 10 and Table 4 summarize our results.

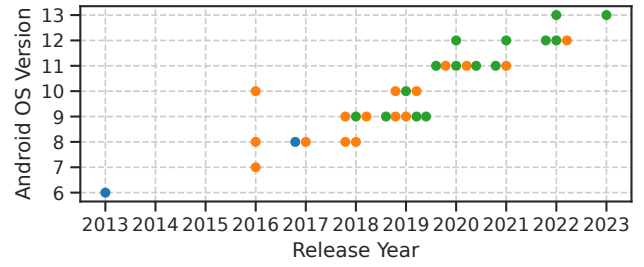


Figure 10: Swarm plot of all surveyed Android devices. We use green ● for StrongBox support, orange ● for TEE support, and blue ● for no support for hardware-backed keys at all.

Device Model	OS	Release	TEE	SE
ASUS Nexus 7 - 2nd Gen (WiFi)	6	2013		
Google Pixel	7	2016	✓	
Google Pixel 2	8	2017	✓	
Google Pixel 3	9	2018	✓	✓
Google Pixel 4 (Unlocked)	10	2019	✓	✓
Google Pixel 4a	11	2020	✓	✓
Google Pixel 5 (Unlocked)	12	2020	✓	✓
Google Pixel 6 (Unlocked)	12	2021	✓	✓
Google Pixel 7	13	2022	✓	✓
LG Stylo 5	9	2019	✓	
OnePlus 8T	11	2020	✓	
Samsung A51	10	2019	✓	
Samsung Galaxy A10s	10	2019	✓	
Samsung Galaxy A13 5G	11	2021	✓	
Samsung Galaxy A40	9	2019	✓	
Samsung Galaxy A7	8	2016	✓	
Samsung Galaxy A71	11	2020	✓	
Samsung Galaxy J7 (2018)	8	2018	✓	
Samsung Galaxy Note 10	9	2019	✓	✓
Samsung Galaxy Note20	11	2020	✓	✓
Samsung Galaxy S10	9	2019	✓	✓
Samsung Galaxy S21 Ultra	11	2021	✓	✓
Samsung Galaxy S22 5G	12	2022	✓	✓
Samsung Galaxy S23	13	2023	✓	✓
Samsung Galaxy S8 (T-Mobile)	8	2017		
Samsung Galaxy S9 (Unlocked)	9	2018	✓	
Samsung Galaxy Tab A 10.1	10	2016	✓	
Samsung Galaxy Tab S4	8	2018	✓	
Samsung Galaxy Tab S6 (WiFi)	9	2019	✓	✓
Samsung Galaxy Tab S7	11	2020	✓	✓
Samsung Galaxy Tab S8	12	2022	✓	✓
Sony Xperia XZ3	9	2018	✓	
Xiaomi 12 Pro	12	2022	✓	

Table 4: All 33 surveyed devices from the AWS device farm. The *Device Model* is the name provided by the AWS API. The *OS Version* refers to the tested OS version and the device might be available with different versions.

These results show that the vast majority of devices released after 2020 provide access to an SE. And all of those have at least TEE support. The earliest device with SE support is the Google Pixel 3 which was released 2018. The release year and Android OS version strongly correlate meaning that newer devices indeed run more recent OS versions. In our sample, all devices with OS version 12+ (Android S, API 31) offer an SE. However, based on the available API distribution (Figure 8), only 14% of active handhelds run this OS version. For an estimate of the overall availability, we take the relative presence of SEs for each OS version and weigh it based on the API distribution. This yields an estimate of 45% devices supporting SE globally for data from January 2023 (up from 39% in August 2022). However, this value will differ between countries and we believe it will continue to improve as more device features, such as contactless payments, rely on SEs.

Note that the percentages are cumulative and in order to get the actual percentage for a specific version number, e.g. API 9 in January 2023, we have to calculate  $81.2\% - 68.0\% = 13.2\%$ . For calculating the estimated share of supported devices we first compute the prevalence of SE support per API level using Table 4. This yields the following data: API 9 (50.0%) API 10 (25.0%) API 11 (57.1%) API 12 (80.0%).

We then calculate for January 2023 as follows:  $13.2\% \cdot 50.0\% + 19.5\% \cdot 25.0\% + 24.4\% \cdot 57.1\% + 24.1\% \cdot 80.0\% \approx 44.7\%$ . And for August 2022 likewise:  $14.5\% \cdot 50.0\% + 22.3\% \cdot 25.0\% + 27.0\% \cdot 57.1\% + 13.5\% \cdot 80.0\% \approx 39.1\%$ .

## B LIMITATIONS OF THE IOS API

The iOS API for using the Secure Enclave allows for P256 private-key pairs and a limited set of operations, namely key agreement (provide a public key and receive the ECDH result) and signing (provide a digest and receive a signature)<sup>8</sup>. There is no support for symmetric keys or operations as with StrongBox on Android. Nevertheless, the results from the asymmetric operations can be used with other cryptographic algorithms. For instance, the Apple API supports using a handle to a key stored inside the SE together with ECIES encryption algorithms like `eciesEncryptionCofactorVariableIVX963SHA256AESGCM`. We use this for the HiddenSloth implementation on iOS.

However, using such ECIES encryption for emulating a scheme similar to LongSloth is not possible due to two obstacles. (i) The first concern is that the documentation is inconclusive<sup>9</sup> regarding whether all operations (including the AEC-GCM algorithm) are performed inside the SE. For safety, and fearing that this may differ across devices and releases, we assume that this is not the case. If the AES-GCM algorithm is performed outside the SE, then an adversary could capture its state after the first block and compute more blocks on an external computer. Hence, we cannot increase

the difficult by requiring a certain length for the input or output stream.

(ii) Even if we assume that the full ECIES scheme is performed inside the SE, we cannot construct a MAC/KDF using either the encryption or decryption operations. The encrypt operation uses a random ephemeral key. Therefore, the resulting ciphertext will be different every time and thus we cannot use it for deriving keys in a LongSloth-way. The decrypt operation verifies the authenticity (the GCM tag) of the input. Therefore, it fails for arbitrary input.

We believe that we therefore have to rely on combining multiple ECDH operations as we do in our RainbowSloth construction. For the HiddenSloth ratchet operations the two listed concerns are non-issues since (i) we perform the re-encryption before adversary control and (ii) we want to perform authenticated encryption.

## C SECURITY PROOFS

We prove the security of LongSloth (Section 4.1), RainbowSloth (Section 4.2), and HiddenSloth (Section 5.2).

### C.1 Security proofs for LongSloth

We first proof LongSloth Hardness as we will reuse the result for our proof of LongSloth Indistinguishability.

#### C.1.1 LongSloth Hardness.

**PROOF.** We constructively proof Theorem 3 by deriving the success rate of  $\mathcal{A}$  with a given wall time budget  $B$ . For this we first note that  $k$  only allows  $\mathcal{A}$  to verify its guesses, but provides no helpful information otherwise. This is because the pre-image resistance of HKDF implies  $\mathcal{A}$  does not learn any information about  $\omega_{post}$  (and thus any of the previous state) from  $k$ . Hence, knowing  $k$  provides  $\mathcal{A}$  with no advantage when choosing  $pw'$  candidates.

Second, we show that  $\mathcal{A}$  has to consider each  $pw'$  candidate independently. For this we observe that the output  $\omega_{pre} = \text{PwHASH}(\pi.salt, pw')$  is indistinguishable from random as per our assumptions, i.e. the input to `SE.HMAC` is independent for each  $pw$ . We note that the salt sampled as per  $\Xi$ . `KEYGEN` rules out any pre-computations by  $\mathcal{A}$ .

Third, we show that  $\mathcal{A}$  has to pay the full costs ( $l \cdot c_{\text{HMAC}}$ ) for each of their guesses. This follows from the *existential unforgeability under adaptive chosen-messages* [25, p.113] that we can assume for `SE.HMAC`. In particular,  $\mathcal{A}$  does not gain any information from submitting a prefix of  $\omega_{pre}$ . With the pre-image resistance of HKDF, this implies that  $\mathcal{A}$  must perform all steps of the  $\Xi$ . `DERIVE`( $\psi, pw', h$ ) using oracle  $O_{SE}$ . Hence, one password guess reduces ( $l \cdot c_{\text{HMAC}}$ ) units from the adversary's budget  $B$ .

Fourth, we derive the probability of success for a set password guesses. Let  $X$  be a random variable drawn from the distribution  $\mathcal{P}$  by the challenger to determine the password  $pw$ . Assume to the advantage of  $\mathcal{A}$  that they can efficiently sample a set  $G$  of passwords from  $\mathcal{P}$  with each  $g_i \in G$  independently drawn from  $\mathcal{P}$  as random variable  $Y$  with  $g_i \neq g_j \forall i, j$ . Then the probability that one of the guesses allows  $\mathcal{A}$  to derive  $k$  is:  $\Pr[\mathcal{A} \text{ wins}] = \sum_{pw \in \mathcal{P}} \sum_{g \in G} \Pr[X = pw] \cdot \Pr[Y = pw'] = \sum_{pw \in \mathcal{P}} \Pr[X = pw] \cdot \sum_{g \in G} \Pr[Y = pw'] = |G| \cdot \Pr[Y = pw'] \leq \frac{|G|}{2^m}$ . Based on our arguments above, the maximum size of  $G$  that the adversary can

<sup>8</sup>In particular, the Secure Enclave: [...] Works only with NIST P-256 elliptic curve keys. These keys can only be used for creating and verifying cryptographic signatures, or for elliptic curve Diffie-Hellman key exchange (and by extension, symmetric encryption)". See: [https://developer.apple.com/documentation/security/certificate\\_key\\_and\\_trust\\_services/keys/protecting\\_keys\\_with\\_the\\_secure\\_enclave](https://developer.apple.com/documentation/security/certificate_key_and_trust_services/keys/protecting_keys_with_the_secure_enclave)

<sup>9</sup>The API design hints that AES is performed outside the SE: The `SecureEnclave.P256.KeyAgreement` type allows access to a `x963DerivedSymmetricKey` from outside the SE. See: [https://developer.apple.com/documentation/cryptokit/sharedsecret/x963derivedsymmetrickey\(using:sharedinfo:outputbytecount:\)](https://developer.apple.com/documentation/cryptokit/sharedsecret/x963derivedsymmetrickey(using:sharedinfo:outputbytecount:))



verify given budget  $B$  is  $|G| = \frac{B}{(l \cdot \text{CHMAC})}$ . Substituting in the previous equation yields:  $\Pr[\text{KeyHard}_{A,\Xi} = 1] \leq \frac{B}{(l \cdot \text{CHMAC})} \cdot \frac{1}{2^m}$ .  $\square$

### C.1.2 LongSloth Indistinguishability.

**PROOF.** We prove Theorem 1 by reduction using the IND-HMAC Experiment. However, the classic definitions assume a high-entropy key  $k$  sampled uniformly at random from  $\{0, 1\}^\lambda$ . In the LongSloth algorithm, this key is derived from a password sampled from  $\mathcal{P}$ . Hence, we separately account for the adversary's using their wall-time budget to guess  $pw$  and hence derive  $k$  using their Oracle access. The success chance is captured by the Hardness result  $\frac{B}{\sigma \cdot 2^m}$  from above.

The LongSloth protocol  $\Xi$  runs on a SE with HMAC support SE-WITH-HMAC (Definition 3); let SE-OP be the HMAC protocol run by the SE. For the reduction, let's assume there exists an efficient adversary  $\mathcal{A}_\Xi$  against  $\Xi$ ; we build an adversary  $\mathcal{A}$  against SE-OP.  $\mathcal{A}$  interacts with a SE-OP-challenger  $C$  and simulates a  $\Xi$ -challenger to  $\mathcal{A}_\Xi$ .

Definition 16 recalls the IND-HMAC $_{\mathcal{A}}$  experiment played by  $\mathcal{A}$  and  $C$ . Bellare proved that HMAC is a PRF under the sole assumption that its underlying compression function is a PRF [12]. As in the original paper by Bellare, it is convenient to consider a PRF-adversary  $\mathcal{A}$  that takes inputs (Section 3.2 of [12]). Figure 11 illustrates how to leverage  $\mathcal{A}_\Xi$  to build an efficient adversary  $\mathcal{A}$  breaking the IND-HMAC security of SE-OP (Definition 17).

**Definition 16** (IND-HMAC Experiment [12]). Let  $\lambda$  be a fixed security parameter. Let  $\mathcal{A}$  be a WT adversary and  $C$  a WT challenger. The SE-OP indistinguishability experiment IND-HMAC $_{\mathcal{A}}$  is defined as follows:

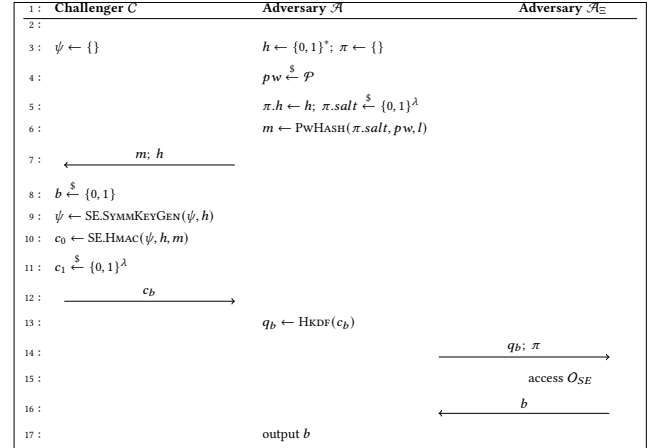
- (1) Let the state  $\psi$  be freshly initialized and  $h$  an arbitrary (but fixed) key handle.  $C$  randomly samples a secret key  $k \leftarrow \text{SE.HMACKeyGen}(\psi, h)$  and sets  $\psi.h \leftarrow h$ .
- (2)  $\mathcal{A}$  receives oracle access SE.HMAC under the WT conditions (Definition 2).
- (3)  $\mathcal{A}$  submits a chosen plaintext  $m \in M$  to  $C$ .
- (4)  $C$  randomly samples  $b \xleftarrow{\$} \{0, 1\}$ , and provides  $\mathcal{A}$  with  $c_0 \leftarrow \text{SE.HMAC}(\psi, h, m)$  if  $b = 0$  or  $c_1 \xleftarrow{\$} \{0, 1\}^\lambda$  otherwise.
- (5)  $\mathcal{A}$  outputs a bit  $b'$  and wins iff  $b = b'$ .
- (6) The experiment returns 1 iff  $\mathcal{A}$  wins, otherwise 0.

**Definition 17** (IND-HMAC Security [12]). A SE-OP protocol is IND-HMAC secure if for all WT adversaries  $\mathcal{A}$  with time budget  $B$ , there is a function  $\text{NEGL}$  such that for all  $\lambda$  and an HMAC operation with length  $l$  having a cost of  $\sigma(l)^{10}$ ,

$$\Pr[\text{IND-HMAC}_{\mathcal{A}} = 1] \leq \frac{1}{2} + \text{NEGL}(\lambda) + \frac{B}{\sigma(l) \cdot 2^m}$$

The challenger  $C$  starts by initializing its internal state  $\psi$  with a secret key. Adversary  $\mathcal{A}$  selects a  $m$ -entropy secure password  $pw$  and uses it to generate a message for the challenger  $C$ ; it sets  $m \leftarrow \text{PwHASH}(\text{salt}, pw, l)$  (where  $\text{salt}$  is a random salt). Challenger  $C$  samples a random bit. If  $b = 0$  it provides  $\mathcal{A}$  with the HMAC of  $m$ , i.e.,  $c_0 \leftarrow \text{HMAC}(\psi, h, m)$ ; otherwise it samples a random bit string

<sup>10</sup>In the main part of the paper we assume a fixed  $n$  (or  $l$  for RainbowSloth) chosen based on the application requirements. In our proofs we parameterize  $\sigma$  to simplify the reduction.



**Figure 11: LongSloth security reduction.**  $\mathcal{A}$  leverages the efficient adversary  $\mathcal{A}_\Xi$  to play against  $C$  and break the IND-HMAC security of SE-OP.

$c_1 \xleftarrow{\$} \{0, 1\}^\lambda$  of the same size as the HMAC output. In order to guess whether  $c_b$  is the output of a HMAC or a random bit string, the adversary  $\mathcal{A}$  provides  $\mathcal{A}_\Xi$  with the stretched key  $q_b \leftarrow \text{HKDF}(c_b)$ .  $\mathcal{A}_\Xi$  determines whether  $q_b$  originated from a key stretching output or a random source. To this purpose,  $\mathcal{A}_\Xi$  can access the oracle  $O_{SE}$  under the WT conditions (see Definition 2). It finally returns  $b = 1$  if it believes  $q_b$  originated from a random source and  $b = 0$  otherwise. Finally, the adversary  $\mathcal{A}$  deduces that  $C$  computed the HMAC of  $m$  if  $b = 0$  and that it sampled a random bit string if  $b = 1$ .

We observe that  $\mathcal{A}$  wins the IND-HMAC $_{\mathcal{A}}$  experiment with the same probability as  $\mathcal{A}_\Xi$  wins  $\text{KEYIND}_{\mathcal{A}_\Xi}$ . As a result, the existence of an efficient adversary  $\mathcal{A}_\Xi$  winning  $\text{KEYIND}_{\mathcal{A}_\Xi}$  with probability  $p > 1/2 + \text{NEGL}(\lambda) + \frac{B}{\sigma(l) \cdot 2^m}$  implies the existence of an efficient adversary  $\mathcal{A}$  winning IND-HMAC $_{\mathcal{A}}$  with the same probability  $p$ . This directly violates the assumption that SE-OP runs a IND-HMAC secure HMAC algorithm, hence a contradiction.  $\square$

## C.2 Security proofs for RainbowSloth

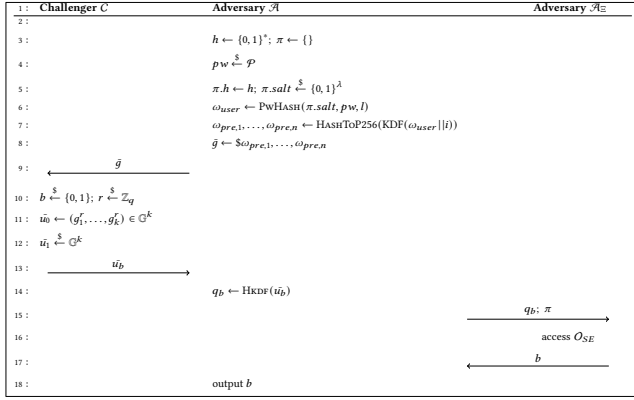
### C.2.1 RainbowSloth Hardness.

**PROOF.** The proof for RainbowSloth is analogous to the one for LongSloth (see §C.1.1) with the cost of the critical operation exchanged for  $n \cdot c_{\text{ECDH}}$  as per Definition 4.  $\square$

### C.2.2 RainbowSloth Indistinguishability.

**PROOF.** We prove Theorem 2 by reduction using the Generalized DDH Experiment. However, the classic definitions assume a high-entropy key  $k$  sampled uniformly at random from  $\{0, 1\}^\lambda$ . In the LongSloth algorithm, this key is derived from a password sampled from  $\mathcal{P}$ . Hence, we separately account for the adversary's using their wall-time budget to guess  $pw$  and hence derive  $k$  using their Oracle access. The success chance is captured by the Hardness result  $\frac{B}{\sigma \cdot 2^m}$  from above.

The RainbowSloth protocol  $\Xi$  runs on a SE with ECDH support SE-WITH-ECDH (Definition 4); let SE-OP be the DDH secure



**Figure 12: RainbowSloth security reduction.**  $\mathcal{A}$  leverages the efficient adversary  $\mathcal{A}_\Xi$  to play against  $C$  and break the generalized DDH assumption.

Diffie-Hellman key exchange run by the SE. For the reduction, let's assume there exists an efficient adversary  $\mathcal{A}_\Xi$  against  $\Xi$ ; we build an adversary  $\mathcal{A}$  against SE-OP.  $\mathcal{A}$  interacts with a SE-OP-challenger  $C$  and simulates a  $\Xi$ -challenger to  $\mathcal{A}_\Xi$ .

Definition 18 recalls the (generalized) decisional Diffie-Hellman (DDH) [11] experiment played by  $\mathcal{A}$  and  $C$ . Figure 12 illustrates how to leverage  $\mathcal{A}_\Xi$  to build an adversary  $\mathcal{A}$  breaking the DDH assumption (Definition 19).

**Definition 18** (Generalized DDH [11]). Let  $\lambda$  be a fixed security parameter. Let  $n$  be an integer and  $\mathbb{G}$  a large cyclic group of prime order  $q$ . Let  $\mathcal{A}$  be a WT adversary and  $C$  a WT challenger. The generalized DDH indistinguishability experiment  $\text{DDH}_{\mathcal{A}}$  is defined as follows:

- (1)  $\mathcal{A}$  provides  $C$  with  $(g_1, \dots, g_n) \in \mathbb{G}^n$ .
- (2)  $C$  randomly samples  $b \xleftarrow{\$} \{0, 1\}$ . If  $b = 0$  it randomly samples  $r \xleftarrow{\$} \mathbb{Z}_q$  and sets  $\tilde{u}_0 = (g_1^r, \dots, g_n^r) \in \mathbb{G}^n$ ; otherwise it randomly samples  $\tilde{u}_1 = (u_1, \dots, u_n) \in \mathbb{G}^n$ . It then provides  $\mathcal{A}$  with  $\tilde{u}_b$ .
- (3)  $\mathcal{A}$  outputs a bit  $b'$  and wins iff  $b = b'$ .
- (4) The experiment returns 1 iff  $\mathcal{A}$  wins, otherwise 0.

**Definition 19** (Generalized DDH [11]). The generalized decisional Diffie-Hellman assumption holds in  $\mathbb{G}$  if for any  $n$ , there is a function  $\text{NEGL}$  and  $n$  ECDH operations having cost of  $\sigma(n)$ , such that

$$\Pr[\text{DDH}_{\mathcal{A}} = 1] \leq \frac{1}{2} + \text{NEGL}(\lambda) + \frac{B}{\sigma(n) \cdot 2^m}$$

Figure 12 shows how to leverage the RainbowSloth adversary  $\mathcal{A}_\Xi$  to build an efficient adversary breaking the the generalized DDH assumption (Definition 19).

The adversary  $\mathcal{A}$  selects a  $m$ -entropy secure password  $pw$ . It then converts it into a  $n$  P-256 public keys as described in Algorithm 2: it picks a random salt  $salt \xleftarrow{\$} \{0, 1\}^*$ , computes  $\omega_{user} \leftarrow \text{PwHASH}(salt, pw, l)$  and  $\omega_{pre,i} \leftarrow \text{HASHToP256}(\text{KDF}(\omega_{user} || i))$  for  $i \in [0, \dots, n]$ . It then sends  $\tilde{g} = (\omega_{pre,0}, \dots, \omega_{pre,n})$  to the challenger  $C$ .

$C$  samples a random bit  $b$  and a field element  $r \xleftarrow{\$} \mathbb{Z}_q$ . If  $b = 0$ , it computes  $\tilde{u}_0 = (g_1^r, \dots, g_n^r)$ ; otherwise it randomly samples  $\tilde{u}_1 \in \mathbb{G}^n$ . It finally sends  $\tilde{u}_b$  to  $\mathcal{A}$ .

In order to guess the bit  $b$  picked by the challenger,  $\mathcal{A}$  provides  $\mathcal{A}_\Xi$  with the stretched key  $q = \text{KDF}(\tilde{u}_b)$ .

$\mathcal{A}_\Xi$  determines whether  $q_b$  originated from its password  $pw$  or a random source (leveraging its access to the  $O_{SE}$  oracle). It returns  $b = 0$  if it believes  $q_b$  originated from  $pw$  and  $b = 1$  otherwise. Finally, the adversary  $\mathcal{A}$  deduces that  $C$  computed  $u_0$  from  $\tilde{g}$  if  $b = 0$  and randomly generated  $\tilde{u}_1$  if  $b = 1$ .

We observe that  $\mathcal{A}$  wins the  $\text{DDH}_{\mathcal{A}}$  experiment with the same probability as  $\mathcal{A}_\Xi$  wins  $\text{KEYIND}_{\mathcal{A}_\Xi}$ . As a result, the existence of an efficient adversary  $\mathcal{A}_\Xi$  winning  $\text{KEYIND}_{\mathcal{A}_\Xi}$  with probability  $p > 1/2 + \text{NEGL} + \frac{B}{\sigma(n) \cdot 2^m}$  implies the existence of an efficient adversary  $\mathcal{A}$  winning  $\text{DDH}_{\mathcal{A}}$  with the same probability  $p$ . This directly violates the assumption that the key exchange SE-OP is secure under DDH, hence a contradiction.  $\square$

### C.3 Security proof of HiddenSloth

Both 1S-HiddenSloth and MS-HiddenSloth run on top of a key stretching scheme such as LongSloth (Section 4.1) and RainbowSloth (Section 4.2); as a result, they run on a SE with either HMAC or ECDH support. As mentioned in Section 5.1 and Section 5.2, any HiddenSloth protocol  $\Delta$  requires an authenticated IND-CPA secure stream cipher AE running in the user space (i.e., outside of the SE).

#### C.3.1 MS-HiddenSloth Indistinguishability.

**PROOF.** We prove Theorem 5 by reduction. Let's assume there exists an efficient adversary  $\mathcal{A}_\Delta$  against  $\Delta$ ; we build an adversary  $\mathcal{A}$  against AE.  $\mathcal{A}$  interacts with a AE-challenger  $C$  and simulates a  $\Delta$ -challenger to  $\mathcal{A}_\Delta$ . Definition 20 recalls the IND-CPA experiment played by  $\mathcal{A}$  and  $C$ ; Definition 21 recalls the definition of AE's IND-CPA security taking into account the ability of the adversary to break the underlying key stretching scheme.

**Definition 20** (AE IND-CPA Experiment). Let  $\lambda$  be a fixed security parameter. Let  $\mathcal{A}$  be a WT adversary with wall time budget  $B$  and  $C$  a WT challenger. Let  $\Xi$  be a  $\sigma$ -hard key stretching scheme. The AE indistinguishability experiment  $\text{IND-CPA}_{\mathcal{A}}$  is defined as follows:

- (1)  $C$  randomly samples a secret key  $k \leftarrow \Xi.\text{KEYGEN}(\psi, pw, h)$ .
- (2)  $\mathcal{A}$  receives oracle access  $\text{AE.ENC}$  under the WT conditions (Definition 2).
- (3)  $\mathcal{A}$  submits two chosen plaintexts  $(m_0, m_1) \in M^2$  to  $C$ .
- (4)  $C$  randomly samples  $b \xleftarrow{\$} \{0, 1\}$  and  $iv \xleftarrow{\$} IV$ , and provides  $\mathcal{A}$  with  $c_b, t_b \leftarrow \text{AE.ENC}(k, iv, m_b)$ .
- (5)  $\mathcal{A}$  receives oracle access  $\text{AE.ENC}$  under the WT conditions.
- (6)  $\mathcal{A}$  outputs a bit  $b'$  and wins iff  $b = b'$ .
- (7) The experiment returns 1 iff  $\mathcal{A}$  wins, otherwise 0.

**Definition 21** (AE IND-CPA). A stream cipher AE keyed with the output of an  $m$ -entropy secure key stretching scheme is IND-CPA secure if for all WT adversaries  $\mathcal{A}$  with time budget  $B$ , there is a function  $\text{NEGL}$  such that for all  $\lambda$ ,

$$\Pr[\text{IND-CPA}_{\mathcal{A}} = 1] \leq \frac{1}{2} + \text{NEGL}(\lambda) + \frac{B}{\sigma \cdot 2^m}$$

It is convenient to grant  $\mathcal{A}$  access to a SE oracle allowing to perfectly simulate SE symmetric encryption operations  $\text{SE.SymmEnc}$

**Algorithm 6** Oracle simulating  $SE.SymmEnc$  operations as performed by a SE-WITH-SYMMENC.

```

1: procedure O.INIT()
2:    $\psi \leftarrow \{\}$ 
3: procedure O.ENCRYPT( $h', m$ )
4:    $\psi \leftarrow SE.SYMMKEYGEN(\psi, h')$ 
5:    $iv \xleftarrow{\$} IV$ 
6:   return  $SE.SYMMENC(\psi, h', iv, m)$ 

```

as performed by a SE with symmetric encryption support SE-WITH-SYMMENC (Definition 5). This oracle is described in Algorithm 6 and is initialized calling O.INIT before starting the experiment.

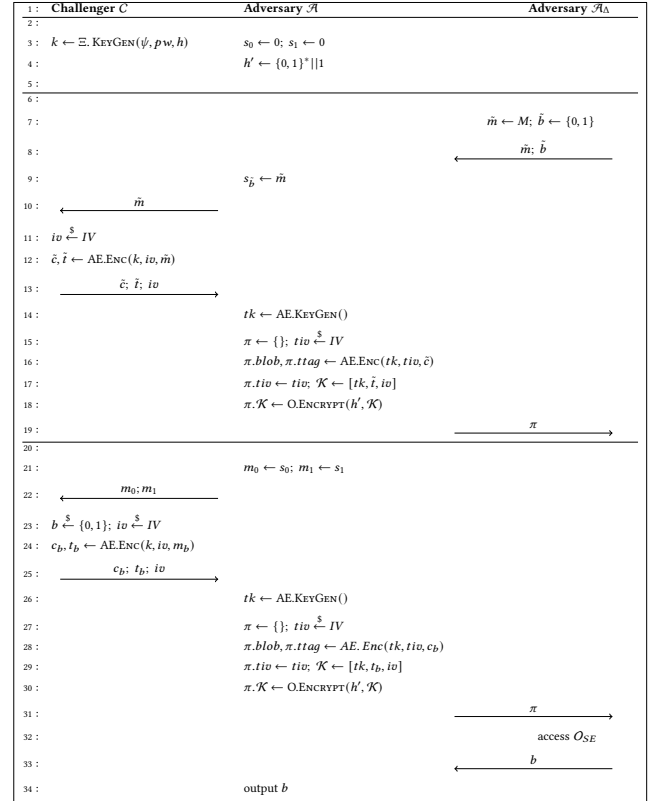
Figure 13 illustrates how to leverage the MS-HiddenSloth adversary  $\mathcal{A}_\Delta$  to build an efficient adversary  $\mathcal{A}$  breaking the IND-CPA security of AE (Definition 21). Figure 13 generates a random key  $k$  calling AE.KEYGEN. The experiment should ideally generate  $k$  through a Sloth key stretching algorithm but Theorems 1 and 2 state that  $\mathcal{A}_\Delta$  cannot distinguish a key generated by a Sloth algorithm from a purely random key.

The key insights of the proof are the following. (i) Challenger  $C$  indirectly simulates Algorithm 4 for adversary  $\mathcal{A}$ . (ii) Adversary  $\mathcal{A}$  does not need to run a decryption operation before Line 26 as it previously cached the latest messages  $\tilde{m}$  from  $\mathcal{A}_\Delta$  within its variables  $s_0$  and  $s_1$ . It can thus feed them to  $C$  during the attack phase and only needs to simulate the outer encryption layer (Line 29 to Line 34 of Algorithm 5).

During the preparation phase (lines 7 to 21 of Figure 13), the adversary  $\mathcal{A}_\Delta$  provides  $\mathcal{A}$  with a message  $\tilde{m} \leftarrow M$  and a bit  $\tilde{b}$ .  $\mathcal{A}$  makes use of its oracle access to AE.ENC (step 2 of Definition 20) by forwarding  $\tilde{m}$  to the challenger  $C$ .  $C$  encrypts  $\tilde{m}$  and replies with the corresponding  $iv$ , ciphertext  $\tilde{c}$ , and tag  $\tilde{t}$ .  $\mathcal{A}$  stores the queried message as  $s_b \leftarrow \tilde{m}$ ; it will later use it to simulate the state of the protocol  $\Delta$  to  $\mathcal{A}_\Delta$ . It then sample a fresh key  $tk$  to re-encrypt  $\tilde{c}$ , which simulates the outer encryption layer of Algorithm 5.  $\mathcal{A}$  makes use of its SE encryption oracle (defined in Algorithm 6) to encrypt  $tk$ ,  $\tilde{t}$ , and  $iv$ .  $\mathcal{A}$  now holds all information to craft a state  $\pi$  as if generated by a  $\Delta$ -challenger. The preparation phase repeats a number of times chosen by  $\mathcal{A}_\Delta$  (under WT conditions).

$\mathcal{A}$  eventually prepares two messages for the challenger  $C$ :  $m_0 \leftarrow s_0$  and  $m_1 \leftarrow s_1$ . The values  $s_0$  and  $s_1$  retain the latest messages queried by  $\mathcal{A}_\Delta$  during the preparation phase for  $\tilde{b} = 0$  and  $\tilde{b} = 1$ , respectively. The challenger  $C$  encrypts either  $m_0$  or  $m_1$  based on a random bit  $b$  and provides  $\mathcal{A}$  with the corresponding ciphertext and tag  $c_b, t_b \leftarrow AE.Enc(k, iv, m_b)$  and  $iv$ . At this point,  $c_b$  is the encryption of the latest pair  $(\tilde{m}, b)$  queried by  $\mathcal{A}_\Delta$  during the preparation phase (at Line 7).  $\mathcal{A}$  finally re-encrypt those information to simulate the outer encryption layer similarly to the preparation phase.

$\mathcal{A}_\Delta$  accesses the oracle  $O_{SE}$  and eventually determines whether  $\pi$  contains the encryption of the latest pair  $(\tilde{m}_0, 0)$  or  $(\tilde{m}_1, 1)$  it submitted during the preparation phase. It eventually returns  $b = 0$  if it believes  $\pi$  contains the encryption of  $\tilde{m}_0$  and  $b = 1$  otherwise. Finally, the adversary  $\mathcal{A}$  deduces that  $C$  encrypted  $m_0$  if  $b = 0$  and  $m_1$  if  $b = 1$ .



**Figure 13: MS-HiddenSloth security reduction.**  $\mathcal{A}$  leverages the efficient adversary  $\mathcal{A}_\Delta$  to play against  $C$  and break the IND-CPA security of AE.

We observe that  $\mathcal{A}$  wins the IND-CPA $_{\mathcal{A}}$  experiment with the same probability as  $\mathcal{A}_\Delta$  wins experiment DE-MS-IND $_{\mathcal{A}_\Delta}$ . As a result, the existence of an efficient adversary  $\mathcal{A}_\Delta$  winning experiment DE-MS-IND $_{\mathcal{A}_\Delta}$  with probability  $p > 1/2 + \text{NEGL}(\lambda) + \frac{B}{\sigma \cdot 2^m}$  implies the existence of an efficient adversary  $\mathcal{A}$  winning IND-CPA $_{\mathcal{A}}$  with the same probability  $p$ . This directly violates the assumption that AE is IND-CPA secure, hence a contradiction.  $\square$

### C.3.2 MS-HiddenSloth Hardness.

**PROOF.** We prove Theorem 6 by showing that the probability that adversary  $\mathcal{A}$  wins the deniable encryption hardness experiment for protocol  $\Delta$  is no bigger than its probability of winning the hardness experiment against its underlying key stretching scheme  $\Xi$ . That is,  $\Pr[\text{KeyHard}_{\mathcal{A}, \Delta} = 1] \leq \Pr[\text{KeyHard}_{\mathcal{A}, \Xi} = 1]$ .

First,  $\mathcal{A}$  removes the outer encryption layer calling  $SE.SYMMDEC$  and  $AE.DEC$  (Line 27 and Line 28 of Algorithm 5). As a result,  $\pi.blob$  contains the encryption of  $data$ , encrypted using a key derived from  $pw$ . This operation requires a one-time cost of  $T_{SE.SYMMENC}$  from the adversary's budget  $B$  and does not need to be repeated for each guess.

Second, we observe that  $\pi$  (and  $\pi.blob$  in particular) only allows  $\mathcal{A}$  to verify their guesses, but provides no helpful information otherwise. This observation follows from the IND-CPA resistance of

the underlying AE encryption scheme. Hence, knowing  $\pi$  provides  $\mathcal{A}$  with no advantage when choosing  $pw'$  candidates.

Third, we show that  $\mathcal{A}$  has to pay the full cost  $c_{\Xi}$  required to access the SE and run the underling key stretching scheme  $\Xi$  used by  $\Delta$ . DECRYPT.  $\mathcal{A}$  must find a key  $k$  such that  $\text{AE.DEC}(k, \pi.iv, \pi.blob, \pi.tag)$  returns *data* (Line 17 of Algorithm 4). Assuming AE is a permutation-based cipher, there is a single  $k$  satisfying this property. Let's assume (to the advantage of the adversary) that  $\mathcal{A}$  can brute-force AE.DEC to recover  $k$  (since it does not require access to the SE and thus does not cost  $\mathcal{A}$ 's budget).  $\mathcal{A}$  must now win the key stretching hardness experiment presented in Definition 9 against a  $\Xi$ -challenger. Assuming  $\Xi$  is  $\sigma_{\Xi}$ -hard, Definition 10 indicates  $\mathcal{A}$  pays budget  $c_{\Xi} = \sigma_{\Xi}$  for each of their guesses.

Finally, since the adversary pays a one-time cost  $T_{\text{SE}, \text{SYMMENC}}$  and a cost of  $\sigma_{\Xi}$  for each of their guesses, the overall cost of guessing  $pw$  is  $\sigma > \sigma_{\Xi}$  per guess. It follows that  $\frac{B}{\sigma \cdot 2^m} < \frac{B}{\sigma_{\Xi} \cdot 2^m}$ , and thus  $\Pr[\text{DEHARD}_{\mathcal{A}, \Delta} = 1] \leq \Pr[\text{KEYHARD}_{\mathcal{A}, \Xi} = 1]$ .  $\square$