



# Incrementally Updateable Honey Password Vaults

Haibo Cheng, Wenting Li, and Ping Wang, *Peking University*; Chao-Hsien Chu, *Pennsylvania State University*; Kaitai Liang, *Delft University of Technology*

<https://www.usenix.org/conference/usenixsecurity21/presentation/cheng-haibo>

This paper is included in the Proceedings of the  
30th USENIX Security Symposium.

August 11–13, 2021

978-1-939133-24-3

Open access to the Proceedings of the  
30th USENIX Security Symposium  
is sponsored by USENIX.

# Incrementally Updateable Honey Password Vaults\*

Haibo Cheng<sup>1</sup>, Wenting Li<sup>1</sup>, Ping Wang<sup>1,†</sup>, Chao-Hsien Chu<sup>2</sup>, Kaitai Liang<sup>3</sup>  
<sup>1</sup>Peking University   <sup>2</sup>Pennsylvania State University   <sup>3</sup>Delft University of Technology

## Abstract

Password vault applications allow a user to store multiple passwords in a vault and choose a master password to encrypt the vault. In practice, attackers may steal the storage file of the vault and further compromise all stored passwords by offline guessing the master password. *Honey vaults* have been proposed to address the threat. By producing plausible-looking decoy vaults for wrong master passwords, honey vaults force attackers to shift offline guessing to online verifications.

However, the existing honey vault schemes all suffer from intersection attacks in the multi-leakage case where an old version of the storage file (e.g., a backup) is stolen along with the current version. The attacker can offline identify the decoys and completely break the schemes. We design a generic construction based on a *multi-similar-password model* and further propose an *incremental update* mechanism. With our mechanism, the attacker cannot get any extra advantages from the old storage, and therefore degenerates to an attacker only with knowledge of the current version.

To further evaluate the security in the traditional single-leakage case where only the current version is stolen, we investigate the theoretically optimal strategy for online verifications, and propose practical attacks. Targeting the existing schemes, our attacks crack 33%–55% of real vaults via *only one-time* online guess and achieve 85%–94% accuracy in distinguishing real vaults from decoys. In contrast, our design reduces the values of the two metrics to 2% and 58% (close to the ideal values 0% and 50%), respectively. This indicates that the attackers needs to carry out 2.8x–7.5x online verifications to break our scheme.

## 1 Introduction

*Password vaults*, a.k.a *wallets* or *managers*, are highly recommended for password management. A user can store multiple

\*Due to the page limit, the mathematical derivation and some design details are left in the full version of this paper (see the authors' websites).

†Corresponding author. He is also with Key Laboratory of High Confidence Software Technologies (PKU), Ministry of Education, China.

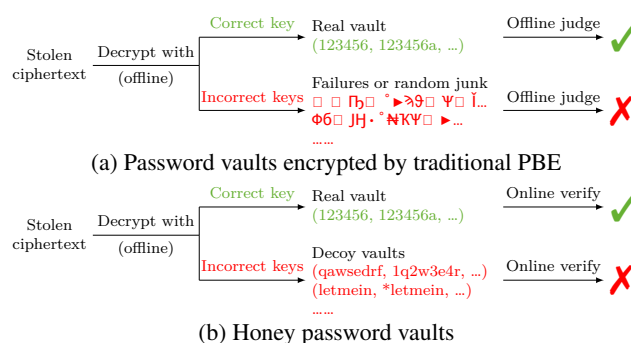


Figure 1: The difference between traditional and honey password vaults in the view of attackers.

passwords in a vault and further set a master password to encrypt the vault. The user thus only needs to remember the master password instead of a long list of daily-use passwords. In practice, the user usually uses the vault among multiple clients (e.g., smartphone or PC), and requires its synchronization via online services. The synchronization may be provided by the vault applications (e.g., LastPass and 1Password) or third-party file sync services (e.g., Dropbox and iCloud). However, the sync services may suffer from leakage [25, 30, 42, 43, 45], which leads to a great threat for password vaults.

If an attacker steals the storage file of a vault (including the ciphertext), the attacker can launch guessing attacks against the master password to compromise all stored passwords. For a vault encrypted by traditional password-based encryption (PBE), decrypting it with an incorrect guess will yield a failure (i.e.,  $\perp$ ) or random junk. So the attacker can immediately identify the validity of guesses *offline*. In addition, since the master password is human-memorable, it may be low-entropy [10, 50] and could be guessed as easily as a website login password [34, 48, 49]. Accordingly, the attacker can efficiently carry out this offline attack with a high probability of success.

*Honey password vault* [8] is proposed to address this threat. Its core idea is to *generate plausible-looking decoy vaults for incorrect guesses to confuse attackers*. As shown in Fig. 1,

launching offline guessing produces many decoy vaults (with a real one), which need to be online verified (i.e., trying to log in with passwords in the vaults). By pushing attackers to online verification, the honey vault mechanism significantly enhances the security of vaults, as the verification can be practically detected and prevented [16, 20, 40].

The design of decoy vaults originates from Bojinov et al. [8]. Their proposed *Kamouflage* pre-generates a *static amount* (e.g., 1,000) of decoy vaults with corresponding decoy master passwords, and further stores them with the real ones. Later, Chatterjee et al. [12] introduced a honey vault scheme *NoCrack* based on *Honey Encryption (HE)* [22]. HE is used to turn a vault to a random-looking bit string called *seed* with a probabilistic encoder—*distribution transforming encoder*—and further encrypt the seed. Due to the “honey” feature provided by HE, using an *arbitrary* wrong master password in decryption can yield a random-looking seed that will be further decoded to a decoy vault on the fly. This brings attackers much more difficulties to tell the real vault, compared with the pre-generating method. Subsequently, Golla et al. [17] proposed *adaptive encoders* which adjust themselves according to the encrypted vault to make decoys more similar to it. Cheng et al. [13] found both Chatterjee et al.’s [12] and Golla et al.’s [17] encoders suffer from *encoding attacks*. They further proposed a generic transformation that can convert a probability model to an encoder resisting encoding attacks.

However, all existing honey vault schemes suffer from *intersection attacks* in the multi-leakage case where an old version of the storage file is stolen along with the current version. This is an open question left in [12, 17]. More specifically, the schemes only provide full update for a vault, i.e., reprocessing the updated vault as a brand new one, even if a user just changes a password (or add a new one). This yields a totally different new version for each decoy vault (by decrypting the new ciphertext with the same master password); and meanwhile the old and new versions of the real vault are the same except for the changed password. Hence, the attacker can offline identify the real vault. This is a *realistic* threat because: 1) the old version of the storage file usually is backed up and stored with the current version by the online services or applications, (for example, Dropbox keeps all history versions of files for 30 days [15], and 1Password automatically creates a backup for each change [6]); 2) the online storage may suffer from multiple leakages due to increasing number of network attacks and software bugs [5, 19, 24, 32, 42–44].

## 1.1 Our Contributions

To resist intersection attacks, we propose a generic construction and an *incremental update* mechanism for HE-based honey vaults. We build our construction from: 1) a *multi-similar-password model* which models the conditional password distribution given multiple old passwords (i.e., the old vault); and 2) the corresponding conditional encoder which

can encode a password given multiple old ones. With the construction, we can encode the changed (or added) password to a (sub) seed and pad it to the tail of the vault seed. With a prefix-keeping PBE scheme, the similarity between the old and new versions of each decoy vault is kept the same as that for the real vault. This means that our scheme resists intersection attacks. Formally, the attacker cannot get any extra advantages from the old storage file, and degenerates to an attacker in the single-leakage case (where the attacker only steals the current version).

To further evaluate the security of (HE-based) honey vault schemes against distinguishing attacks in the single-leakage case, we formally investigate the optimal strategy for online verifications and further propose several practical attacks<sup>1</sup>. For the existing schemes [12, 13, 17], our attacks crack 33%–55% of real vaults via *only one-time* online guess and achieve 85%–94% accuracy in distinguishing real vaults from decoys. We further find that the adaptive encoder [17] does leak extra information about the real vault. Leveraging the information, our attacks can achieve 91%–93% distinguishing accuracy against the adaptive encoder, which is 6.2%–9.0% higher than that of the static variant. This makes the adaptive encoder more insecure than its static variant.

To instantiate our construction, we design a multi-similar-password model according to users’ password-generating habits. The design is built on the top of a single-password model (capturing how the user creates a brand new password), a single-similar-password model (capturing how a user creates a new password by reusing an old one) and an unreused probability function. For our design, the existing and our proposed attacks only crack at most 2% of real vaults via one-time online guess and achieve at most 58% distinguishing accuracy. The results are close to 0% and 50% maintained by an ideal secure scheme, respectively. This also means the attacker has to carry out 2.8x–7.5x online verifications against our scheme as compared to others. Since online verifications can be quickly detected and prevented, our design achieves a significant improvement on security.

In summary, we describe our main contributions as follows.

1. We propose a new generic construction and an incremental update mechanism for HE-based honey vault schemes, which resists intersection attacks.
2. We formally investigate the optimal strategy for online verifications and further propose several practical attacks, which can effectively distinguish real and decoy vaults for the existing honey vault schemes.
3. We instantiate our construction with a well-designed multi-similar-password model, which can generate more plausible-looking decoys.

<sup>1</sup>Here, we only focus on human-generated passwords, since it is of great challenge to generate indistinguishable decoys [12] for them but trivial for randomly-generated ones. Although many vault applications recommend users to use randomly-generated passwords, they always store some human-generated passwords in the vaults [33, 38].

## 2 Background and Related Work

### 2.1 Traditional Solutions to Offline Guessing

A straightforward solution to master password guessing is to leverage a special password hashing as the key derivation function (KDF) used in password-based encryption (PBE), such as an iterated hash function [23,41], a memory-hard function [9,39]. LastPass employs this solution (using the 100,100 rounds of PBKDF2-SHA256 [3]) to increase the computational cost of attackers in launching master password guess. Nevertheless, the cost of a valid user is also increased by the same factor. Without leveraging heavy hashing, one may use an extra key stored on a device (e.g., iOS keychain [4], a server [27]), to enhance the master password to a cryptographic key for further encryption, like 1Password [4]. But this has a defect that if the device gets lost without any backup, all the passwords stored in the vault cannot be recovered anymore. Note that these solutions can be used in honey vault schemes to achieve complementary protection. There may be other approaches but we don't explore them here. We will only focus on the solutions based on honey vaults.

### 2.2 Honey Encryption

*Honey Encryption (HE)* [22], proposed by Juels and Ristenpart, can resist the brute-force attack by yielding plausible-looking messages for arbitrary incorrect keys, even in the case where a low-entropy key (e.g., password) is used. Later, Jaeger et al. [21] proved that HE satisfies the stronger notions of target-distribution semantic security and target-distribution non-malleability. The core design of HE relies on an encoder, called *distribution transforming encoder (DTE)*, being able to capture the message distribution. Intaking a message  $M$  following some distribution  $\mathcal{M}$ , DTE can encode it to a bit string  $S$  called *seed* which is indistinguishable from a random string. HE further encrypts  $S$  to a ciphertext  $C$  by a traditional but carefully-chosen PBE scheme with a key  $K$ . Decrypting  $C$  with a wrong key  $K'$  (e.g., a guessing key from attackers), the carefully-chosen PBE (e.g., the CTR-mode AES with PBKDF used in [12,17,22]) can yield a random-looking bit string  $S'$ . DTE then decodes  $S'$  into a honey message  $M'$  which is sampled from the same distribution  $\mathcal{M}$ . Note in the context of honey vault, the message and key correspond to the password vault and master password, respectively.

Juels and Ristenpart use *inverse sampling* to convert a distribution to an encoder called *IS-DTE*. This method performs well for simple distributions, e.g., uniform distributions. But when handling messages (e.g., natural language) with a huge space size and a complex distribution, it definitely yields explosive complexity in time and storage space. To tackle this problem, Chatterjee et al. [12] introduce a *natural language encoder*, and later Cheng et al. [13] propose a *probability model transforming encoder* (see Sections 2.3 and 2.4).

Table 1: The storage format of honey vaults

Plaintext part	Domain	Facebook	Myspace	000Webhost	Twitter	...
	Username	Aaron	Aaron1	AaronJ	Aaron	...
	Randomly-generated	No	No	Yes	No	...
	Password position	1	3	1	2	...
Ciphertext part	Human-generated	(123456, 123456a, 1234567, ...)				
	Randomly-generated	(cYp97@v84G\$9GNv%3R, ...)				

Note: Each randomly-generated password is encoded by the encoder for the uniform distribution and further encrypted separately. All human-generated passwords are encoded by the encoder for the vault model and further encrypted as a whole.

### 2.3 HE-based Honey Vault Schemes

Unlike traditional solutions to offline master password guessing, honey vault schemes yield decoy vaults for incorrect guesses and therefore force attackers to online verify these decoy vaults. We only focus on HE-based schemes [12,13,17] in this paper due to their advantage on security.

**Storage format.** The existing schemes only use HE to encrypt passwords and leave other parts (e.g., domains and usernames) in plaintext<sup>2</sup>. We give an example in Table 1 to show the storage format.

**Vault model.** The probability model for password vaults (vault model, for short) is the foundation used to generate indistinguishable decoys. It should characterize the real vault distribution as precisely as possible. Because of users' various password generation (and reuse) habits, it is a great challenge to design models for human-generated passwords<sup>3</sup>. The existing vault models [12,17] choose a single-password model as the base to characterize the single-password distribution and further extend it to capture the similarity (reuse habits) among multiple passwords in a vault: Chatterjee et al. [12] use the probabilistic context-free grammar (PCFG) model and extend it by the sub-grammar approach; Golla et al. [17] leverage the Markov model and extend it by the reuse-rate approach. Note that Cheng et al. [13] do not propose a new vault model but leverage/recommend Golla et al.'s design. In addition, Golla et al. introduce and apply *adaptive* concept to vault model, encoder and honey vault scheme. Before encrypting a real vault  $V$ , an adaptive scheme adjusts its vault model (as well as its encoder) according to  $V$ . With well-designed adjustments, an adaptive model may produce decoys that are more similar to  $V$ , bringing more difficulties to attackers in identifying them.

<sup>2</sup>One may choose to further encrypt domains and usernames, but this will break the "honey" property of HE. Note that a real username is registered on the domain, but its decoy is not. Thus, the attacker can easily identify the decoy vaults by registering with the usernames. Without using encryption, Chatterjee et al. [12] provide another way to hide domains. They generate decoy accounts for the domains where users have not registered, and further store the decoy usernames in plaintext. But this method still suffers from the aforementioned attack. It seems that there does not exist an effective solution to hide domains and usernames. We leave this as an open problem.

<sup>3</sup>Since it is trivial for randomly-generated passwords [12], we do not consider them in this paper.

**Encoder.** An encoder for a vault model should encode a vault sampled from the model to a seed being indistinguishable from a random bit string. However, the natural language encoders designed by Chatterjee et al. [12] and used by Golla et al. [17] fail to achieve this requirement, and therefore suffer from encoding attacks [13]. Cheng et al. [13] tackle this vulnerability by employing their encoders for the old models. To evaluate the existing honey vault schemes without the negative effect caused by encoding attacks, we will adopt Cheng et al.’s encoders to [12] and [17] in this paper, and still refer to the resulting schemes as Chatterjee et al.’s and Golla et al.’s. Under our adoption, Golla et al.’s scheme (with Cheng et al.’s encoder) becomes the same as Cheng et al.’s scheme (with Golla et al.’s model as they recommended).

**Deployment consideration.** Due to the special feature of honey vaults, if a user enters an incorrect master password (e.g., a typo), it will get a decoy vault and further lead to a login failure. Dynamic security skin can be used to address the issue as suggested in [8, 12]. This approach shows a picture to the user according to the master password input. By checking if the picture is identical to the one from the last (correct) input, the user can verify the correctness of the master password. Unlike the ciphertext, the picture is not stored by the application, and thus will not be stolen from the online storage. Note the user does not need to remember the whole picture but just a vague impression. Other typo-correcting methods (e.g., [11]) can also be used in deployment without putting an extra burden on users.

## 2.4 Model-to-encoder transformation

Cheng et al. [13] propose a generic method to transform an arbitrary probability model to a *probability model transforming encoder*, which resists encoding attacks. Their core idea is to assume that messages are created by generating paths (i.e., a sequence of generating rules). Based on the idea, they formalize all current models for the single password or password vault. For example, in their formalization for PCFG models, the generating rules are production rules and the generating paths are leftmost derivations. To further encode a message, their encoder parses *all* generating paths of the message, randomly selects one path with its probability, and encodes each rule in the path. (In contrast, the existing encoders in [12, 17] use deterministic path selection, therefore it is easy to exclude the decoy seeds of which paths are not the deterministic ones.) In this way, each seed of this message can be randomly and uniformly picked. This feature is called *seed uniformity*, which is the “cure” to encoding attacks.

However, in some models associated with great ambiguity, a message may be generated by various paths. Parsing all these paths may yield heavy time complexity in encoding. Although Cheng et al. attempted to reduce the ambiguity by pruning some unnecessary paths, the low encoding performance still limits the scalability of their encoders.

## 3 Our Incrementally Updateable Scheme

We propose a generic construction for vault models and further construct an encoder, which provides incremental update for password vaults and achieves the update security (i.e., resisting intersection attacks).

### 3.1 Our New Construction

In practice, the passwords in a vault  $V = (pw_i)_{i=1}^n$  are generated one by one. Therefore, the probability  $\Pr_{\text{real}}(V)$  can be expanded as

$$\Pr_{\text{real}}(V) = \prod_{i=0}^{n-1} \Pr_{\text{real}}(pw_{i+1} \mid pw_1, pw_2, \dots, pw_i), \quad (1)$$

where  $\Pr_{\text{real}}(pw_{i+1} \mid (pw_{i'})_{i'=1}^i)$  is the probability of creating a new password  $pw_{i+1}$  under the condition of given  $i$  old passwords  $(pw_{i'})_{i'=1}^i$ . Naturally, we can leverage a conditional probability model  $\Pr_{\text{MSPM}}(\cdot \mid \cdot)$  to estimate the conditional probability and further construct a vault model. We denote  $\Pr_{\text{MSPM}}(\cdot \mid \cdot)$  as *multi-similar password model*.

### 3.2 Conditional Probability Model Transforming Encoder

For a probability model with the following construction

$$\Pr_{\text{model}}((M_i)_{i=1}^n) = \prod_{i=1}^n \Pr_{\text{model}}(M_i \mid (M_{i'})_{i'=1}^{i-1}), \quad (2)$$

using Cheng et al.’s model-to-encoder transformation can yield a probability model transforming encoder. However, the encoder has exponential time complexity if the model is ambiguous. Specifically, if there exist  $k_i$  paths to generate  $M_i$  from  $(M_{i'})_{i'=1}^{i-1}$ , then there will be  $\prod_{i=1}^n k_i$  generating paths for  $M = (M_i)_{i=1}^n$ . Cheng et al.’s encoder has to calculate the probabilities of  $\prod_{i=1}^n k_i$  paths and randomly select one with its probability, which yields the time complexity  $O(\prod_{i=1}^n k_i)$ .

To reduce the time complexity of the encoder, we extend Cheng et al.’s [13] transformation for conditional probability models. Since a conditional probability model can be seen as a probability model for each condition, a *conditional probability model transforming encoder (conditional encoder for short)* can be achieved by transforming the conditional probability model for each condition with Cheng et al.’s transformation.

By using the conditional encoder ( $\text{encode}(\cdot \mid \cdot)$ ,  $\text{decode}(\cdot \mid \cdot)$ ) for  $\Pr_{\text{model}}(\cdot \mid \cdot)$ , we design a new encoder for  $\Pr_{\text{model}}(\cdot)$  as follows:

1. To encode a message  $M = (M_i)_{i=1}^n$ : encode  $M_i$  to a seed  $S_i$  by  $\text{encode}(M_i \mid (M_{i'})_{i'=1}^{i-1})$ , then output the concatenating of  $\{S_i\}_{i=1}^n$ , i.e.,  $S = S_1 \parallel S_2 \parallel \dots \parallel S_n$ .
2. To decode a seed  $S$ : split  $S$  to  $\{S_i\}_{i=1}^n$  according to the fixed length of  $S_i$ , decode  $M_i$  from  $S_i$  in order by  $\text{decode}(S_i \mid (M_{i'})_{i'=1}^{i-1})$ , then output  $M = (M_i)_{i=1}^n$ .

In contrast to Cheng et al.’s encoder, our new encoder combined by conditional encoder only needs to select one path from  $k_i$  paths for  $1 \leq i \leq n$ , which significantly reduces the time complexity to  $O(\sum_{i=1}^n k_i)$ . In addition, since the conditional encoder is seed-uniform, our new encoder naturally inherits this property and therefore resists encoding attacks.

### 3.3 Incrementally Updateable Encoder

The proposed conditional encoder for the multi-similar-password model can encode a vault password by password<sup>4</sup>. This naturally brings an incremental update mechanism. In detail, we update the storage file of a vault as follows:

1. To add a new password: decrypt the vault, encode the new password by the conditional encoder for the multi-similar-password model, add the password seed to the tail of the vault seed, record the password position in plaintext, and finally encrypt the updated seed with the same (correct) key and the same nonce.
2. To delete an old password: mark the password as deleted (in plaintext) without changing the ciphertext.
3. To change an old password: delete the old password, add the new password as in Item 1, and update the password position for the corresponding account.

To achieve the update security, the PBE adopted in HE must satisfy the *prefix-keeping* property:

1. If a string  $str_1$  is a prefix of a string  $str_2$ , then the ciphertext  $C_1$  of  $str_1$  is also a prefix of the ciphertext  $C_2$  of  $str_2$  with the same key and the same random nonce. Note the nonce is stored in plaintext, e.g., the salt for PBKDF and the initialization vector for CTR-mode.
2. If a ciphertext  $C_1$  is a prefix of a ciphertext  $C_2$  (with the same nonce), then the plaintext  $str_1$  of  $C_1$  is a prefix of the plaintext  $str_2$  of  $C_2$  under any (incorrect) key.

The PBE scheme used in [12, 13, 17], i.e., AES in CTR-mode merged with PBKDF, satisfies the prefix-keeping property. Thus, we use the same scheme for our design.

**Intersection-attack resistance.** Compared with the full update in the current honey vault schemes, our design decreases the time complexity in updating but also guarantees the update security against intersection attacks. As shown in Table 2, owing to the prefix-keeping property that our chosen PBE provides, we ensure that the old ciphertext is a prefix of its new version. Decrypting the old and new ciphertexts with the same (incorrect) master password, we will have that the new seed is the same as its old version except for the added tail (note if one leverages the schemes with full update, these two seeds will be totally different). Accordingly, the new decoy vault is identical to its old version except for the changed or added passwords. This means that the similarity of the old and new versions for each decoy vault is kept the same as that for the real vault. Therefore, if an attacker steals the old

<sup>4</sup>When initializing a vault with a set of existing passwords, the application can shuffle the passwords and then encode them one by one.

Table 2: The difference between old and new vaults after changing Facebook password<sup>1</sup>

	Position <sup>2</sup>	Ciphertext	Seed <sup>3</sup>	Passwords <sup>3</sup>
Previous	1	$C_1    \dots    C_{10}$	$S_1    \dots    S_{10}$	$pw_1, \dots, pw_{10}$
Updated	11	$C_1    \dots    C_{10}    C_{11}$	$S_1    \dots    S_{10}    S_{11}$	$pw_1, \dots, pw_{10}, pw_{11}$

<sup>1</sup> We take the vault in Table 1 as an example and assume it contains 10 (human-generated) passwords.

<sup>2</sup> By position we refer to the password position of Facebook account.

<sup>3</sup> The previous and updated seeds/passwords are under a master password which may be the correct one or an incorrect one.

version of the vault storage file along with the current version, the attacker cannot get extra advantages and degenerates to an attacker only with the knowledge of the current version. This means the attacker cannot launch intersection attacks but only (traditional) distinguishing attacks based on the current version. We confirm this statement in an experiment with real-world datasets (see Appendix C).

**Other potential threats.** Unlike the existing schemes, our scheme maintains the password history (including the changed and deleted passwords), which may bring an advantage in distinguishing attacks. Based on users’ password-changing habits, the attacker may leverage similarity among the old and new passwords for distinguishing. (Note the password-similarity attack proposed in Section 4 can be naturally extended for this purpose.) To address this issue, we can leverage a well-designed multi-similar-password model to capture the password-changing habits and further generate plausible-looking password histories for decoys. In addition, keeping the password history is provided by many real-world vault applications as a feature (rather than a flaw), e.g., Last-Pass [2]. Our scheme naturally provides this feature, while the existing ones cannot.

### 3.4 Multi-Similar-Password Model

To instantiate our construction for honey vault scheme, we need a *multi-similar-password model*  $\text{Pr}_{\text{MSPM}}$  which can precisely estimate  $\text{Pr}_{\text{real}}(pw_{i+1} | (pw_{i'})_{i'=1}^i)$ . To the best of our knowledge, such models do not exist in the literature. Pal et al. [36] mention this notion in their work on password guessing, but they leave it as future work without providing a specific model. Here, we give a simple design for the model.

**Our design.** We model a user’s generation of a new password  $pw_{i+1}$  with  $i$  old passwords  $(pw_{i'})_{i'=1}^i$  by a simplification that  $pw_{i+1}$  either is created by “reusing” an old password (including a direct reuse or a slight modification of it) or is a brand new creation. Accordingly,

$$\begin{aligned} & \text{Pr}_{\text{MSPM}}(pw_{i+1} | pw_1, pw_2, \dots, pw_i) \\ &= f(i) \text{Pr}_{\text{SPM}}(pw_{i+1}) + \frac{1-f(i)}{i} \sum_{i'=1}^i \text{Pr}_{\text{SSPM}}(pw_{i+1} | pw_{i'}), \end{aligned} \quad (3)$$

where  $\text{Pr}_{\text{SSPM}}$ ,  $\text{Pr}_{\text{SPM}}$  and  $f(i)$  represent a *single-similar-*

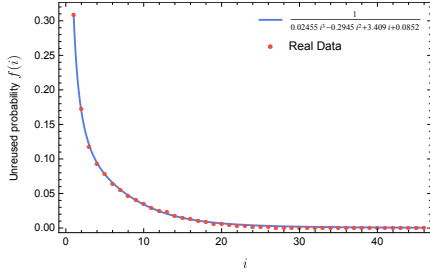


Figure 2: Unreused probability  $f(i)$  that the  $i+1$ -th password is not reused from the first  $i$  passwords.

password model, a single-password model, and the unreused probability function, respectively. The single-password model captures the new generation, the single-similar-password model captures the reusing, and the unreused probability function is the probability that the user does not reuse  $i$  old ones.

Then we carefully instantiate the three components. For the single-password model, we choose to use a Markov model with well-set parameters (denoted as *Best-Markov*), since it performs the best under the single-password attack among existing single-password models [34, 35, 46, 49] (see Appendix D). For the single-similar-password model, we design a simple model which only captures the most prevalent password-reuse habit, i.e., head or tail modification [14]. For simplicity, this model regards two passwords as reused passwords if the length of their longest common substring is at least half of the maximum length of them (note we say the password pair has Feature LCSStr). In this way, our model is simple and further leads to the efficiency of encoding. The details of the model are given in Appendix A. We also try to use the existing single-similar-password models, e.g., the password-to-path model [36] and the context Wasserstein autoencoder [37]. But we find our model is more suitable for honey vaults because of its best performance on the decoy vault generation and the high encoding efficiency (see Appendix E).

For the unreused probability function, we can leverage a real-world vault dataset (Pastebin, see Section 5.2) and count the empirical probability  $\hat{f}(i)$  of the event that the  $i+1$ -th password is not reused from the first  $i$  passwords<sup>5</sup>. Further, we perform nonlinear regression on  $\hat{f}(i)$  and find  $\hat{f}(i)$  can be fitted well with a 3-degree rational function in the form  $f(i) = 1 / (\sum_{k=0}^3 a_k i^k)$ . As shown in Fig. 2, the fitted function  $f_{\text{fit}}(i) = \frac{1}{0.02455i^3 - 0.2945i^2 + 3.409i + 0.0852}$  is very close to  $\hat{f}(i)$ , and  $|f_{\text{fit}}(i) - \hat{f}(i)| \leq 4.101 \times 10^{-3}$ . Thus, we use  $f_{\text{fit}}(i)$  as the unreused probability function<sup>6</sup>. In addition,  $\hat{f}(i)$  decreases as  $i$  increases, indicating that the more passwords a user has, the

<sup>5</sup>There may be a great number of  $i+1$ -tuples for some  $i$ . Thus, we estimate the probability by sampling  $10^4$  tuples (with replacement) and counting  $\hat{f}(i)$  in the samples instead of directly counting in all tuples.

<sup>6</sup>In the experiments with 5-fold cross-validation (see Section 5.2), the coefficients of  $f_{\text{fit}}(i)$  are not obtained from Pastebin but the training set, which is 4/5 of Pastebin.

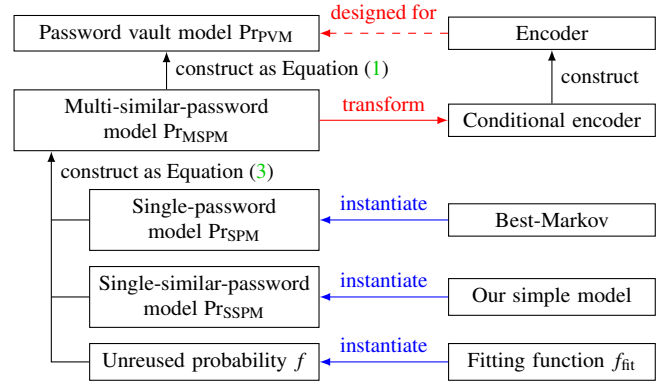


Figure 3: Technical roadmap of our designs.

less likely he is to create a brand new password. This reflects the limit of human memory on passwords.

With the above constructions and instantiations, we finally construct a concrete model and an encoder for honey vaults. The full construction is shown in Fig. 3.

**Time complexity of encoding.** Applying the transformation proposed in Section 3.2, we can get a conditional encoder for the multi-similar-password model and with efficiency in encoding. Specifically,  $pw_i$  is generated by reusing old passwords  $(pw_{i'})_{i'=1}^{i-1}$  or is a brand new creation in our model. Therefore, there are at most  $i$  generating paths for  $pw_i$  under the condition of  $(pw_{i'})_{i'=1}^{i-1}$ . The time complexity of encoding a vault  $(pw_i)_{i=1}^n$  is  $O(n^2)$ . The details of the conditional encoder are provided in Appendices A and B.

**Complying with password policies.** Many websites may adopt password policies to prevent users from using weak passwords, e.g., requiring passwords to contain at least 8 characters. The passwords generated by the vault models (i.e., in the decoy vaults) should always comply with the corresponding policies, otherwise, attackers will easily identify those decoys which do not. We introduce an efficient method to adjust our model to support two classic policies, length restriction (e.g.,  $\geq 8$  characters) and character requirement (e.g., the inclusion of an upper-case letter), guaranteeing that all passwords sampled from the model achieve the complying requirement. The core idea is to exclude the noncompliant lengths or characters when encoding and decoding. Specifically, we adjust the length distribution in best-Markov and our single-similar-password model; and meanwhile, we model the position of the required character type and adjust the corresponding character distribution (by adjusting, we mean excluding the lengths or characters not complying the requirement and re-normalizing the probabilities of the rest ones).

We note it is very challenging to guarantee the need w.r.t. more complex policies (e.g., blacklist or password strength requirement) and we leave this as an open problem. For the websites with these policies, users can use randomly-generated passwords that are easily complied with the policies.

### 3.5 Leakage Detection

We propose a mechanism to detect the leakages of storage files of honey vaults. The core idea is to generate and store some decoy accounts (called *honeypot accounts*) in a user’s real vault. For example, if the user has a real account(name) “Alice07” on Google, the vault application may generate (and register) a honeypot account “Alice07” on Yahoo (with a password generated by our model). These honeypot accounts will not be used by the user (note we can choose the websites the user rarely visits for honeypot accounts to avoid the user’s misuse); and meanwhile they are really registered on the corresponding websites. An attacker with the stolen storage file cannot tell them from real accounts and will probably log in to them (for online verification). Once the logins of honeypot accounts occur, the leakage can be reliably detected. Then the user should change all passwords in the vault to prevent consequent account compromise. In this way, we significantly mitigate the risk of vault file leakages.

Further, we consider the threat in the case where a password in a vault is leaked as well as the vault storage file. In this case, the attacker can offline tell the real vault by checking if the leaked password is in the vault. Although this is an important and practical threat, it is not considered in [12, 13, 17]. Fortunately, leveraging the leakage detection mechanism for honey vaults and the existing alert mechanisms for password breaches (e.g., [1, 18]), we can detect the password and vault leakages, respectively. This enables users to timely change the leaked passwords or vaults. Our solution is a “pre-action” mechanism for the threat (preventing it from happening) rather than a post-action (resisting the attacks after the threat is there).

The design details and security analysis are given in the full version of this paper.

## 4 Attacks Against Honey Vault Schemes

To evaluate the security of honey vault schemes in the single-leakage case, we investigate the theoretically optimal strategy for online verifications, and further propose several practical attacks.

### 4.1 Attacker Model

**Attacker ability.** We consider a significant threat for honey vaults: an attacker steals the (current) storage file of a vault (e.g., from online sync services), and tries to reveal all stored passwords from the file. The attacker also gets the program of the honey vault scheme, including the HE algorithm and the encoder, since the program should be stored along with the storage file to provide handy service to users. So the attacker can try to decrypt the ciphertext with a dictionary of master password guesses. To distinguish the real and decoy vaults decrypted from the ciphertext, the attacker may leverage some

public information, e.g., leaked datasets, website password policies. Note that if the vault scheme uses a public dataset to train its vault model, the attacker can identify the dataset from the encoder and may further use it to launch attacks.

In this section, we do not consider the cases where the attacker additionally gets an old version of the storage file or a website password in the vault. We have addressed these two cases in Sections 3.3 and 3.5, respectively.

**Attack process.** To reveal passwords from an encrypted vault, the attacker should decrypt the ciphertext  $c$  with a dictionary of the master password guesses  $\{mpw_i\}_{i=1}^N$ , and then obtain a (large) group of vaults  $\{V_i\}_{i=1}^N$ , in which at most one of the vaults is real. To check the correctness of the vaults, the attacker needs to log in with the passwords in the vaults (i.e., online verification).

The effectiveness of the attack depends on 1) the offline guessing order of master passwords and 2) the online verification order of the vaults. The former is mainly related to the strength of the master password, while the latter relies on the security of the encoder (i.e. the indistinguishability of real and decoy vaults). We recall that a master password is a human-memorable password and may suffer from general password guessing attacks [34, 36, 48, 49]. We take the same research direction as [12, 13, 17], focusing on the security of encoders.

We assume attackers will test target vaults (via online verification) in a descending order defined by a priority function. Each vault is assigned a priority value so that attackers first online verify those vaults with *greater* values. We denote this priority function as  $p$  with a subscript representing the name abbreviation of the attack. In practice,  $p(V_i)$  is strongly related to the probability that  $V_i$  is the real vault among  $\{V_i\}_{i=1}^N$ .

### 4.2 Theoretically Optimal Strategy

The theoretically optimal strategy of online verification is to verify the vaults in the descending order of conditional probabilities, where the condition is all the information that attackers have known. In the following, we investigate this strategy by analyzing the conditional probabilities.

We denote the random variables of the (real) master password, the (real) vault, the (real) seed and the ciphertext, as  $MPW$ ,  $V$ ,  $S$  and  $C$ , respectively. Let  $MPW = mpw_i$  denote the event that a user’s real master password is identical to  $mpw_i$  (i.e., the user chooses  $mpw_i$  as the master password  $MPW$ ). Similarly, we define  $V = V_i$ ,  $S = S_i$  and  $C = c$ . To keep consistency with notations in Section 4.1, we define  $mpw_i$ ,  $V_i$ ,  $S_i$ ,  $c$  such that  $S_i$  is decrypted from  $c$  with  $mpw_i$  and  $V_i$  is decoded from  $S_i$ . We further denote the real master password distribution  $\Pr(MPW = mpw_i)$  as  $\Pr_{MPW}(mpw_i)$ , the real vault distribution  $\Pr(V = V_i)$  as  $\Pr_{\text{real}}(V_i)$ , the decoy vault distribution as  $\Pr_{\text{decoy}}(V_i)$  (i.e., the probability of getting  $V_i$  by decoding a random seed), and the probability of encoding  $V_i$  to  $S_i$  as  $\Pr_{\text{encode}}(S_i | V_i)$  (i.e.,  $\Pr(S = S_i | V = V_i)$ ), respectively.



For the (static) honey vault schemes, the only information attackers can learn about the real vault  $V$  is its ciphertext  $c$ . (An adaptive scheme will leak extra information about  $V$ , see Section 4.4.) Therefore, *the optimal strategy is to verify  $V_i$  in the descending order of  $\Pr(V = V_i | C = c)$ .*

For simplicity, we use  $\Pr(MPW = mpw_i | C = c)$  to estimate  $\Pr(V = V_i | C = c)$ <sup>7</sup>. In addition, we notice that the existing honey vault schemes [12, 13, 17] require a user to create a completely new and distinct master password so that the master password is independent of all the passwords in the vault. This requirement is due to the limitation of HE [22] that cannot guarantee the security for dependent key and message distributions. Therefore,  $\Pr(MPW = mpw_i, V = V_i) = \Pr_{MPW}(mpw_i) \Pr_{\text{real}}(V_i)$ . According to the Bayesian theorem, we have the following theorems. Note the proofs of the theorems in this section are given in the full version of this paper.

**Theorem 1.** *For an arbitrary encoder,*

$$\begin{aligned} \Pr(MPW = mpw_i | C = c) \\ = k \cdot \Pr_{MPW}(mpw_i) \Pr_{\text{real}}(V_i) \Pr_{\text{encode}}(S_i | V_i), \end{aligned} \quad (4)$$

where  $k$  is a constant independent of  $i$ .

This theorem shows that Cheng et al.’s strong encoding attack [13] is a degenerate case of the optimal strategy by only considering the last factor  $\Pr_{\text{encode}}(S_i | V_i)$  as the priority function. If an encoder is not seed-uniform (i.e., the seeds of a message are not randomly chosen when encoding the message), the real and decoy vaults can be effectively distinguished by merely exploiting the encoder (i.e., calculating  $\Pr_{\text{encode}}(S_i | V_i)$ ) without any knowledge of the master password and password vault distributions (i.e.,  $\Pr_{MPW}(mpw_i)$ ,  $\Pr_{\text{real}}(V_i)$ ).

**Theorem 2.** *If the encoder is seed-uniform, then*

$$\begin{aligned} \Pr(MPW = mpw_i | C = c) \\ = k \cdot \Pr_{MPW}(mpw_i) \frac{\Pr_{\text{real}}(V_i)}{\Pr_{\text{decoy}}(V_i)}, \end{aligned} \quad (5)$$

where  $k$  is a constant independent of  $i$ .

Theorem 2 indicates that if an encoder is seed-uniform, attackers cannot get any information from the encoder except the decoy vault distribution (i.e.,  $\Pr_{\text{decoy}}(V_i)$ ). This analysis confirms that Cheng et al.’s transformation is secure, i.e., their encoder resists encoding attacks. By applying the secure encoders to the existing honey vault schemes [12, 17], we only need to consider how to hold against *distribution difference attacks*. This type of attack is defined in [13], referring to the attacks that only exploit the difference between the real and decoy distributions (i.e.,  $\Pr_{\text{real}}(V_i)$  and  $\Pr_{\text{decoy}}(V_i)$ ).

<sup>7</sup>The vaults  $V_i$  and  $V_j$  under different master passwords  $mpw_i$  and  $mpw_j$  may be the same (i.e.,  $V_i = V_j$ ). But this only happens with a low probability, especially if the vaults are of a large size, because the space of vaults is much larger than that of master passwords.

Table 3: Examples of real-to-decoy probability ratios

	Vault	Password	Reuse feature <sup>†</sup>		Increased $n$ -gram number
Example	(123456,1234567)	123456	0	1	4
Real probability	— <sup>*</sup>	$10^{-2}$	0.8	0.2	$10^{-7}$
Decoy probability	$10^{-4}$	$5 \times 10^{-3}$	0.4	0.6	$10^{-9}$
Ratio	— <sup>*</sup>	2	2	0.333	100

<sup>\*</sup> It is hard to precisely calculate the real probability and the ratio for a vault. So we use some methods to estimate the ratio for attacks.

<sup>†</sup> Each feature used in our classifier is a Boolean (binary variable).

We use Theorem 2 to present a new vision for the security analysis w.r.t. HE and honey vault schemes. For a seed-uniform encoder, if the decoy distribution is the same as the real one, i.e.,  $\Pr_{\text{decoy}} = \Pr_{\text{real}}$ , then  $\Pr(MPW = mpw_i | C = c) = \Pr_{MPW}(mpw_i)$  (in this case,  $k = 1$ ). Therefore, the mutual information of  $C$  and  $MPW$  is  $I(MPW; C) = H(MPW) - H(MPW | C) = 0$ . This means that the ciphertext  $C$  does not leak any information of the key  $MPW$ , which achieves the ideal security of HE and honey vault schemes.

Without considering  $\Pr_{MPW}$  (as discussed in Section 4.1), *the optimal online verification order for vaults  $\{V_i\}_i$  is the descending order of*

$$\frac{\Pr_{\text{real}}(V_i)}{\Pr_{\text{decoy}}(V_i)}, \quad (6)$$

which is denoted as  $p_{\text{OPT}}$ . The basic idea behind this *real-to-decoy probability ratio*  $p_{\text{OPT}}$  is simple and intuitive. A high  $p_{\text{OPT}}(V_i)$  means the vault model used in the honey vault scheme significantly underestimates the real probability of  $V_i$ . In other words,  $V_i$  is less likely generated by the vault model than being generated by the user. Therefore,  $V_i$  is more likely to be real among  $\{V_i\}_{i=1}^N$ .

### 4.3 Practical Attacks

The optimal online verification with the priority function  $p_{\text{OPT}}$  is hard to be carried out, since it is difficult to precisely calculate the real probability  $\Pr_{\text{real}}(V_i)$  for attackers. This difficulty also hinders the direct use of existing techniques, e.g., Bayesian updating, in calculating  $p_{\text{OPT}}$ .

Leveraging an advanced model seems to be a straightforward method to estimate  $\Pr_{\text{real}}$ . Unfortunately, all current models have defects, which lead to the misestimation of the probabilities [47]. For example, the PCFG model [49] underestimates the passwords with relative components, e.g., “1q2w3e”, because the model assumes these components are independent and does not further consider their relationships [47]. The misestimation of a model will lead to the misestimation of  $p_{\text{OPT}}(V)$  and further significantly decrease the effectiveness of attacks. Note we have tried to use different single-password models to estimate the real single-password distribution, e.g., Markov models [34], but could

not obtain stable and reasonable effectiveness in attacking all other single-password models.

To overcome the difficulty, we use several methods to estimate the real-to-decoy probability ratio  $p_{OPT}$  and further propose several practical attacks as follows. The characterization of both the single-password distribution and the password similarity are the two significant indices of a vault model. Accordingly, our estimations will focus on these two indices.

**Single-password attack.** To capture the difference between real and decoy vaults on the single-password distribution, we use the real-to-decoy probability ratio on the single password to estimate the ratio  $p_{OPT}(V_i)$  on vault. Formally assume that the passwords in a vault are independent (ignoring their similarity),  $p_{OPT}(V_i)$  can be simplified/estimated as

$$\prod_{pw \in V_i} \frac{\Pr_{\text{real}}(pw)}{\Pr_{\text{decoy}}(pw)}, \quad (7)$$

where  $\Pr_{\text{real}}(pw)$  and  $\Pr_{\text{decoy}}(pw)$  represent the real and decoy single-password distributions, respectively.  $\Pr_{\text{decoy}}(pw)$  can be calculated by the single-password model in the targeted honey vault scheme, but we still need to estimate  $\Pr_{\text{real}}(pw)$ .

To estimate  $\Pr_{\text{real}}(pw)$ , we directly use the relative frequency of  $pw$  in a password training set, instead of using some password models. We note this is because we do not want to bring the misestimation of password probability yielded by single-password models to our attacks (as discussed above). According to the law of large numbers, the relative frequency of an event converges (almost surely) to its probability, as the number of experiments approaches infinity. To avoid misestimation incurred by inappropriate training sets, we choose the dataset which has been used to train the single-password model by the honey vault schemes (see Section 5.2).

In addition, smoothing is further needed since some passwords not appearing in the training set have zero frequency. With a carefully-designed smoothing method, we propose an estimation  $p_{SP}(pw)$  for  $\frac{\Pr_{\text{real}}(pw)}{\Pr_{\text{decoy}}(pw)}$  as

$$p_{SP}(pw) = \begin{cases} 1 & \text{if } f_a(pw) \leq f_d \text{ and } \frac{f'_r(pw)}{\Pr_{\text{decoy}}(pw)} > 1, \\ \frac{f'_r(pw)}{\Pr_{\text{decoy}}(pw)} & \text{otherwise,} \end{cases} \quad (8)$$

where  $f_a(pw)$  is the absolute frequency of  $pw$  in the training set,  $n$  is the size of the training set,  $\alpha_s$  is a smoothing parameter,  $f_d$  is a parameter representing the demarcation line between high-frequency and low-frequency passwords, and  $f'_r(pw) = \frac{f_a(pw) + \alpha_s}{n + \alpha_s}$ . Our estimation is similar to maximum likelihood estimation (MLE) with Laplace smoothing. Unlike Laplace smoothing used in [34], our smoothing only adds  $\alpha_s$  for the calculated  $pw$  instead of all passwords in the password space. This is because the password space is extremely large, using Laplace smoothing will make  $f'_r(pw)$  to approach to 0 for all  $pw$ . We note that our smoothing may lead to overestimation for the probabilities of some low-frequency passwords.

Thus, we choose to set  $p_{SP}(pw) = 1$  for passwords with absolute frequency no more than  $f_d$  and  $\frac{f'_r(pw)}{\Pr_{\text{decoy}}(pw)} > 1$ . After a few tries, we finally set  $\alpha_s = 1$  and  $f_d = 5$ .

Using  $p_{SP}(pw)$ , we propose a *single-password attack* with the following priority function

$$p_{SP}(V_i) = \prod_{pw \in V_i} p_{SP}(pw). \quad (9)$$

Informally, this attack gives priority to the vaults of which some passwords are not accurately characterized (their probabilities are underestimated) by the single-password model in the targeted honey vault scheme. According to Theorem 2 and Equation (7), the vaults are more likely to be real.

**Password-similarity attack.** To capture the difference between real and decoy vaults on the password similarity, we use the real-to-decoy probability ratio on some features with a Bernoulli naive Bayes classifier to estimate the ratio  $p_{OPT}(V_i)$ . The features should capture the misestimation of the vault models on the password similarity. With a carefully-chosen feature set  $\mathcal{F}$ ,  $p_{OPT}(V_i)$  can be simplified/estimated as

$$\prod_{F \in \mathcal{F}} \frac{\Pr_{\text{real}}(F = F(V_i))}{\Pr_{\text{decoy}}(F = F(V_i))}, \quad (10)$$

where  $F(V_i)$  is the value of Feature  $F$  for  $V_i$  ( $F(V_i) = 1$  if  $V_i$  has Feature  $F$ , otherwise,  $F(V_i) = 0$ ),  $\Pr_{\text{real}}(F = x)$  is the probability that the value of Feature  $F$  is  $x$  for a real vault, and  $\Pr_{\text{decoy}}(F = x)$  is the probability for a decoy vault.

We demonstrate that the estimation is effective. Unlike  $\Pr_{\text{real}}(V_i)$  which is difficult to be calculated,  $\Pr_{\text{real}}(F = x)$  can be counted from a password vault dataset (counting the proportion of vaults which have Feature  $F$  for  $x = 1$ , and the rest proportion for  $x = 0$ ). The vault dataset (Pastebin, in Section 5.2) we are going to use is small, so we only choose two binary features for the Bayes classifier with four parameters. Similarly,  $\Pr_{\text{decoy}}(F = x)$  can be counted from a set of decoy vaults generated by the encoder (decoding random seeds).

To design appropriate features, we first analyze the characterization of vault models on the password similarity. Recall that a user almost always reuses passwords in different accounts [14], therefore, the passwords in his vault usually are similar. A vault model should precisely capture the similarity and further generate similar (i.e., reused) passwords in decoy vaults. In the existing models, a password pair  $(pw_1, pw_2)$  is treated as similar by simple rules: in Golla et al.'s model [17],  $pw_1, pw_2$  are treated as similar if  $pw_1$  is the same as  $pw_2$  except for the last 5 characters (we say  $(pw_1, pw_2)$  has Feature GM); in Chatterjee et al.'s model [12],  $pw_1, pw_2$  are treated as similar if  $pw_1$  and  $pw_2$  share at least one production rule in their PCFG model (we say  $(pw_1, pw_2)$  has Feature CM). The sample treatment leads to the inaccuracy of the models on password similarity.

To crack a vault model, we define Features M and I:  $(pw_1, pw_2)$  has Feature M if the model treats  $pw_1, pw_2$  as

similar;  $(pw_1, pw_2)$  has Feature I if a user can create  $pw_2$  by reusing  $pw_1$ . Then Features M and I capture the similarity among passwords in decoy vaults and real vaults, respectively. Therefore, the difference between Features M and I can be used to exploit the misestimation of the model. Formally, we define the feature difference as follows.

**Definition 1.** We say  $(pw_1, pw_2)$  has Feature  $AB$ , i.e., the difference of Features  $A$  and  $B$ , if  $(pw_1, pw_2)$  has Feature  $A$  but not Feature  $B$ .

With a well-defined Feature I, we can use  $\mathcal{F} = \{MI, IM\}$  to propose a password-similarity attack. Here, we define that a vault  $V$  has Feature X, if there exist two passwords  $pw_1, pw_2$  in  $V$  such that  $(pw_1, pw_2)$  has Feature X. Note that the model probably overestimates the probability of a vault with Feature MI and underestimates that for a vault with Feature IM.

However, it is difficult to precisely define Feature I. We use Feature LCSStr as an approximation of Feature I to attack Chatterjee et al.’s and Golla et al.’s schemes, due to the fact that modifying the head or tail characters is most popular in reuse habits [14]. But for our scheme, Feature M is Feature LCSStr (see Section 3.4), then Features MI and IM are trivial with Feature LCSStr as Feature I (no vaults have Feature LCSStrLCSStr). So we leverage four password similarity meters used in [14] to define Feature I, including Levenshtein [29], longest common subsequence (LCS), Manhattan [26], and Overlap [28]. For each meter F, we define that  $(pw_1, pw_2)$  has Feature F, if the similarity score of  $(pw_1, pw_2)$  is at least 0.5 under the meter F.

To summarize, we use Equation (10) as the priority function  $p_{PS}$  with  $\mathcal{F} = \{MI, IM\}$  for the *password-similarity attack*. To crack Chatterjee et al.’s scheme, Features M and I are Features CM and LCSStr, respectively; for Golla et al.’s scheme, Features M and I are Features GM and LCSStr, respectively; for our scheme, Features M is Feature LCSStr and Feature I is one of Features Levenshtein, LCS, Manhattan, and Overlap. Note that this attack gives priority to these vaults, in which the password similarity is not well characterized by the vault model in a honey vault scheme. According to Theorem 2 and Equation (10), these vaults are more likely to be real.

**Hybrid attack.** Combining the single-password attack with the password-similarity attack, we propose a hybrid attack with the following priority function

$$p_H(V_i) = p_{SP}(V_i) \cdot p_{PS}(V_i). \quad (11)$$

Note that like  $p_{OPT}$ ,  $p_{SP}$  and  $p_{PS}$  are in the form of real-to-decoy probability ratios, but on different indices. So we keep this form by multiplying  $p_{SP}$  and  $p_{PS}$ , and then the product  $p_H$  can estimate  $p_{OPT}$  more precisely. This is confirmed by our experimental results that the hybrid attack always performs better than the previous two attacks (see Section 5).

**Other attacks.** The support vector machine (SVM) attack and the Kullback–Leibler (KL) divergence attack are proposed by Chatterjee et al. [12] and Golla et al. [17], respectively. Since the latter outperforms the former against all the

existing vault models, we will use the latter for comparison. The priority function of the KL divergence attack is defined as

$$p_{KL}(V_i) = \sum_{j=1}^s f_j \log \frac{f_j}{\Pr_{\text{decoy}}(pw_j)}, \quad (12)$$

where  $V_i$  contains  $s$  unique passwords  $\{pw_j\}_{j=1}^s$ , and  $f_j$  is the relative frequency of  $pw_j$  in  $V_i$ .

Golla et al. [17] exploit extra information to enhance their KL divergence attack, including username, password reuse rate and password policy. Among the three types of information, only password policy has significant improvement for KL divergence attack. Attackers can easily exploit it to distinguish the decoys not complying with the policy. We will also consider the *password policy attack* with a minor difference. Formally speaking, the priority function  $p_{PP}$  of this attack can be defined as: 1) if there exists a password not complying with its policy in the vault  $V_i$ , then  $p_{PP}(V_i) = 0$ ; 2) otherwise,  $p_{PP}(V_i) = 1$ . Unlike [17] which only considers one password in a vault, our password policy attack requires all passwords to comply with their policies and can exclude much more decoy vaults.

#### 4.4 More Attacks to Adaptive Encoders

For static schemes, the storage file  $c$  is the only information that attackers can learn. But for adaptive schemes, attackers can learn extra information about the real vault from the encoder. This is because the adaptive encoder is adjusted according to the encrypted real vault. We exploit the “extra” information and propose more attacks to adaptive schemes.

**Theoretically optimal strategy.** We denote the random adaptive encoder as  $DTE$ , and let  $DTE = DTE^*$  be the event that the encoder is adjusted to  $DTE^*$  according to the real vault. *The theoretically optimal strategy here is to verify the vaults  $\{V_i\}_i$  in the descending order of  $\Pr(V = V_i | DTE = DTE^*, C = c)$ .* Recall that the adaptive encoder  $DTE^*$  and the ciphertext  $c$  are the only information that attackers can learn. Similar to the attacks against static encoders, we leverage  $\Pr(MPW = mpw_i | DTE = DTE^*, C = c)$  to estimate  $\Pr(V = V_i | DTE = DTE^*, C = c)$  and have the following theorem.

**Theorem 3.** *If the adaptive encoder  $DTE^*$  is seed-uniform, then*

$$\Pr(MPW = mpw_i | DTE = DTE^*, C = c) = k \cdot \Pr_{DTE}(DTE^* | V_i) \Pr_{MPW}(mpw_i) \frac{\Pr_{\text{real}}(V_i)}{\Pr_{DTE^*}(V_i)}, \quad (13)$$

where  $k$  is a constant independent of  $i$ , and  $\Pr_{DTE^*}(V_i)$  represents the distribution of decoy vaults generated by  $DTE^*$ ,  $\Pr_{DTE}(DTE^* | V_i)$  represents the probability that  $DTE$  is adjusted to  $DTE^*$  according to  $V_i$ .

The priority functions of optimal strategy for the static and adaptive encoders differ by one factor  $\Pr_{DTE}(DTE^* | V_i)$ , which indicates the “extra” information leaked by  $DTE^*$ .

**Practical attacks.** To carry out practical attacks, we need to calculate  $\Pr_{\text{DTE}}(DTE^* | V_i)$  individually. But this is difficult for Golla et al.’s adaptive encoder [17]. We propose to use a simple method to estimate its value, denote the estimator as  $p_{\text{AE}}(V_i)$ . With  $p_{\text{AE}}$  as the priority function, we propose an *adaptive extra attack*. The basic idea of the estimation is to leverage the real-to-decoy probability ratio on the number of  $n$ -grams whose probability is increased by Golla et al.’s adjustment. We leave the complex mathematical analysis about  $\Pr_{\text{DTE}}(DTE^* | V_i)$  and the details of the estimator  $p_{\text{AE}}(V_i)$  in the full version of this paper. It is worth mentioning that if  $p_{\text{AE}}(V_i) = 0$ , then  $V_i$  must be decoy. As can be seen from this case, exploiting the information leaked by the adaptive encoder, attackers can easily exclude some decoy vaults.

Furthermore, we propose an *adaptive hybrid attack* by combining the adaptive extra attack with the hybrid attack. Its priority function is defined as

$$p_{\text{AH}}(V_i) = \text{sgn}(p_{\text{AE}}(V_i)) \cdot p_{\text{H}}(V_i), \quad (14)$$

where  $\text{sgn}$  is the sign function. At the first attempt, we used  $p_{\text{AE}}(V_i) \cdot p_{\text{H}}(V_i)$ . But we later found out that its performance (sometimes) was worse than that of the hybrid attack. This may be caused by the estimation error. To optimize its performance, we then use  $\text{sgn}(p_{\text{AE}}(V_i))$  for  $p_{\text{AH}}(V_i)$  instead of  $p_{\text{AE}}(V_i)$ . Specifically, the adaptive hybrid attack first excludes the decoy vaults with  $p_{\text{AE}}$  of 0 and then cracks the remaining vaults by launching the hybrid attack. Thus, the adaptive hybrid attack should always outperform the hybrid attack.

## 5 Security Evaluation under Our Attacks

We evaluate the existing and our honey vault schemes over the attacks proposed in Section 4 via real-world datasets. The experimental results show that our scheme achieves a significant improvement on security.

### 5.1 Security Metrics

An attack is more effective if it can use a smaller number of online verifications to identify the real vault for a given ciphertext. This number of online verifications is identical to the rank of the real vault among a large number of decoys in the order defined by the priority function. Thus, we use the *ranks* of real vaults to indicate the security of a honey vault scheme against the attack, as in [12, 13, 17].

Chatterjee et al. [12] and Golla et al. [17] use the *average rank*  $\bar{r}$  as a crucial security metric in their evaluation. Chatterjee et al. [12] also define the *accuracy*  $\alpha$  in distinguishability. Please note that  $\alpha$  is the probability of identifying the real from *only one* decoy (by sorting these two with the priority function), not from a larger number of decoys. To present a comprehensive evaluation, Cheng et al. [13] leverage the *cumulative distribution functions (CDFs)*  $F(x)$  of the ranks.

Note that each incorrect master password yields a decoy. Since the master password space is large, it is difficult to calculate the rank of a real vault by generating all decoys. Instead, we choose Cheng et al.’s method [13] to estimate the rank in relative form by sampling  $N$  decoys ( $N = 999$ ). The rank then is defined as the ratio of the rank to the number of decoys, which is a real number in  $[0, 1]$  and reflects the relative position in the online verification order. For example, a vault of rank 0.2 will be online verified after checking 20% decoys. In relative form (hereafter, by rank we mean its relative form),  $\bar{r}$  and  $\alpha$  can be derived from  $F(x)$  [13] as

$$\bar{r} = 1 - \int_0^1 F(x) dx, \quad \alpha = 1 - \bar{r}. \quad (15)$$

For the sake of comparison fairness, we use the (above) same metrics in our experiments, including  $\bar{r}$ ,  $\alpha$ ,  $F(x)$ . In addition, we also use  $F(0)$  as in [13], which indicates the proportion of real vaults with *rank 0*—the vaults cracked in *only one-time* online verification (i.e., one guess).

A perfectly secure honey vault scheme guarantees that real and decoy vaults should be indistinguishable, so that the ranks under any attacks follow the uniform distribution  $U[0, 1]$  (i.e., any attacks perform the same as the randomly guessing attack with a constant priority function). We then have  $F(x) = F_U(x)$  ( $= x$  for  $0 \leq x \leq 1$ ) and  $\bar{r} = \alpha = 0.5$ . Therefore, we use  $F_U(x)$  as the baseline for the comparison.

### 5.2 Experimental Settings

To present a fair and comprehensive comparison, we utilize the same datasets used in [12, 13, 17]: RockYou as the password dataset and Pastebin as the password vault dataset. RockYou, which is one of the largest leaked plaintext password sets, provides 32.6 million password samples. Being able to maintain the completeness of samples and offer a sufficiently large sample size, it is widely used in the security evaluation on recent password researches [7, 34, 35, 48]. To the best of our knowledge, Pastebin is the only publicly available password vault dataset. It consists of 276 vaults with sizes of 2–50. The data of Pastebin, collected by malware embedded on clients, may indirectly provide us a vision of current exploit means of attackers. In the experiments, we only use these datasets to perform security evaluations. From this perspective, the datasets will bring no harm to valid users and the evaluation results will inspire us to design more secure schemes.

To evaluate the security of honey vault schemes, we do 5-fold cross-validation on Pastebin. Specifically, we randomly divide Pastebin into five parts. We take one part as the test set and the union of other parts as the training set. The vaults in the test set are treated as the real vaults which will be protected (i.e., encrypted) by honey vault schemes and be cracked by attacks. The training set (with RockYou) is used to train the vault models by honey vault schemes. As discussed in Section 4.1, we also use the same training set to train attacks. In this

Table 4: The average rank  $\bar{r}$  of real vaults under attacks

Scheme	KL divergence	Single password	Password similarity	Hybrid
Chatterjee et al.'s [12]	14%	10%	22%	6%
Golla et al.'s [17] (static, $10^0$ )	48%	22%	26%	14%
Golla et al.'s [17] (static, $10^{-1}$ )	37%	19%	23%	14%
Golla et al.'s [17] (static, $10^{-2}$ )	34%	20%	26%	14%
Golla et al.'s [17] (static, $10^{-4}$ )	31%	20%	23%	15%
Golla et al.'s [17] (static, $10^{-6}$ )	30%	19%	24%	14%
Golla et al.'s [17] (static, $10^{-8}$ )	29%	19%	24%	15%
Golla et al.'s [17] (static, $10^{-10}$ )	29%	19%	23%	14%
Golla et al.'s [17] (adaptive, $10^0$ )	54%	22%	25%	14%
Golla et al.'s [17] (adaptive, $10^{-1}$ )	43%	20%	25%	13%
Golla et al.'s [17] (adaptive, $10^{-2}$ )	40%	21%	25%	13%
Golla et al.'s [17] (adaptive, $10^{-4}$ )	37%	20%	25%	13%
Golla et al.'s [17] (adaptive, $10^{-6}$ )	36%	21%	26%	13%
Golla et al.'s [17] (adaptive, $10^{-8}$ )	35%	20%	24%	12%
Golla et al.'s [17] (adaptive, $10^{-10}$ )	34%	20%	24%	13%
Ours	42%	48%	43%	42%

<sup>1</sup>  $10^i$  represents the pseudocount of Laplace smoothing used in Golla-Markov [17].

<sup>2</sup> The average rank  $\bar{r}$  and the accuracy  $\alpha$  have the relationship:  $\bar{r} + \alpha = 1$ .

Table 5: The average rank  $\bar{r}$  under extra attacks against Golla et al. adaptive schemes [17] with different pseudocounts

Attack	$10^0$	$10^{-1}$	$10^{-2}$	$10^{-4}$	$10^{-6}$	$10^{-8}$	$10^{-10}$
Adaptive extra	29%	26%	24%	24%	24%	25%	26%
Adaptive hybrid	9%	8%	8%	8%	7%	7%	7%

setting, we exploit the honey vault schemes to generate decoy and launch attacks to get the rank of each vault in the test set. For each part of Pastebin, we do the above experiment to get ranks of all vaults in Pastebin.

Some important details need to be noticed:

1. The honey vault schemes usually need a password dataset to train their single-password model. As in [12, 17], we adopt RockYou for this purpose. This dataset is also used for attacks (if needed).
2. Some attacks need to calculate decoy probabilities (on single password or password feature). This calculation can be launched with the stolen encoders and does not need an extra dataset.
3. Golla et al. [17] do not provide a training method for their reuse-rate approach. For a fair comparison with previous studies, we directly use their parameters in our experiments without training.

### 5.3 Experimental Results

**The performance of our attacks.** As shown in Fig. 4, Tables 4 and 5, our proposed attacks perform well against all of the existing schemes. For all the static schemes, the hybrid attack has the best performance, achieving 94% accuracy  $\alpha (= 1 - \bar{r})$  against Chatterjee et al.'s scheme [12] and 85%–86% against Golla et al.'s scheme [17] (with different parameters), as the

Table 6: RCDFs  $F(x)$  for honey vault schemes under the corresponding best attacks

Scheme	$F(0)$	$F(1/4)$	$F(1/2)$	$F(3/4)$
Chatterjee et al.'s [12]	55%	93%	98%	99%
Golla et al.'s [17] (static, $10^0$ )	33%	71%	92%	100%
Golla et al.'s [17] (adaptive, $10^0$ )	45%	84%	99%	100%
Ours	2%	37%	61%	80%

$F(0)$  indicates the cracked proportion of vaults via *only one* online guess.

average rank  $\bar{r}$  for Chatterjee et al.'s scheme is 6% and those for Golla et al.'s scheme are 14%–15%. Note Cheng et al.'s honey vault scheme [13] is the same as Golla et al.'s, since we adopt Cheng et al.'s encoder for Golla et al.'s scheme (see Section 2.3). Thus the two schemes achieve the same experimental results, and we do not illustrate the results for Cheng et al.'s scheme separately. With regard to the adaptive scheme, the adaptive hybrid attack outperforms others, capturing 91%–93% accuracy  $\alpha$ , as the average ranks  $\bar{r}$  are 7%–9%. *Our attacks are based on the theoretically optimal strategy with more accurate estimation, yielding stable and high accuracy.*

The performance of the KL divergence attack is severely affected by the pseudocount of Laplace smoothing used in Golla-Markov. While pseudocount is set to 1, we have the worst attack performance, achieving 46% and 52% accuracy against Golla et al.'s static and adaptive encoders [12], respectively. As shown in Figs. 4b and 4c, the RCDFs are close to the baseline, which means the attack performs close to the randomly guessing attack. We demonstrate that *it is difficult to distinguish real vaults from decoys without any information of the real vault distribution*. The attack, however, only estimates the distance between the vault to be sorted and decoy vaults, not considering the distance between the vault and real vaults. This may lead to some misjudgments. There may exist a target with a “large” distance from the decoy vaults but also with a “larger” distance from the real vaults. The attacker will mistreat the target as the real vault which is actually more likely to be a decoy. To propose more effective attacks, we must exploit both the real and decoy distributions.

**The security of the existing schemes.** The accuracy of our proposed attacks reveals the vulnerability of the existing schemes. *In terms of the single-password distribution and the password similarity, the existing schemes fail to characterize the real vault distribution*. This is proved by our experimental results, the single-password and the password-similarity attacks achieving 78%–90% and 74%–78% accuracy, respectively. Further, we find out that the pseudocount has an impact on the security of Golla et al.'s schemes [17]: when it is set to 1 these schemes achieve the best security. Here, we only show the RCDFs with this pseudocount.

Exploiting the extra information leaked by the adaptive scheme, our adaptive hybrid attack increases  $\alpha$  to 91%–93%, which is higher than the  $\alpha$  on the static schemes with the same pseudocounts under the hybrid attack, i.e., 85%–86%. *The*

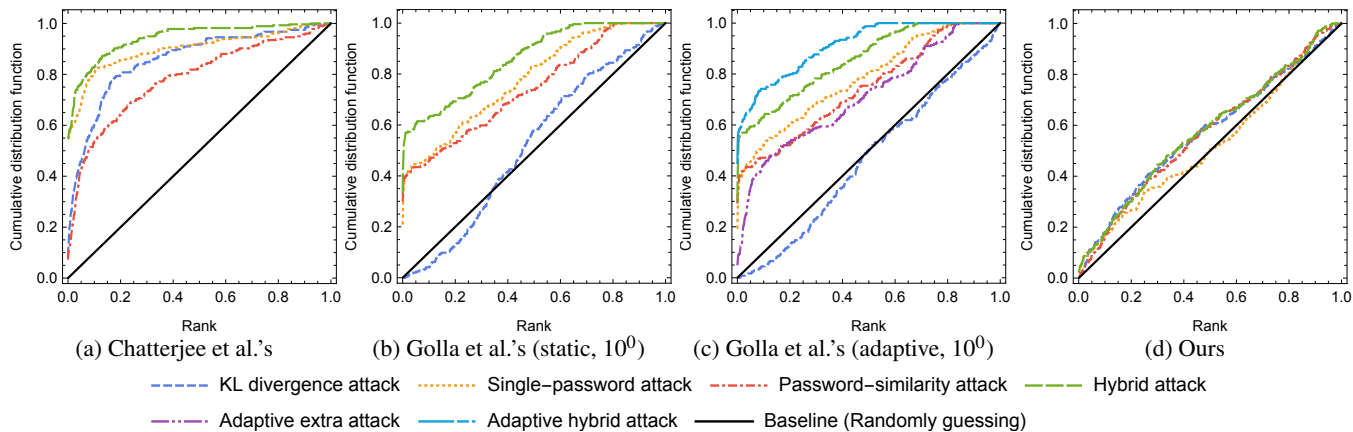


Figure 4: RCDFs for honey vault schemes under attacks.

findings of our experiments overturn the conclusion stated in [12]: “adaptive schemes are more secure than the static ones.” In fact, an adaptive scheme inevitably does leak information of the encrypted real vault. By only exploiting the (leaked) information, our adaptive extra attack achieves 71%–76% accuracy without any knowledge of the real vault distribution.

We note that 1) any adaptive scheme (more precisely, its vault model) is adjusted from its original version according to the real vault (which is about to be encrypted by the scheme); 2) the original model has already been trained with a real password vault dataset. If the size of the training set is sufficiently large, adding one more real vault (i.e., the encrypted vault) cannot significantly improve the precision of the model. As shown in Table 4, the average ranks of static and adaptive schemes using the same pseudocount are almost identical under our hybrid attack (without exploiting the leaked information). We thus conclude that *compared with its static variant, an adaptive scheme cannot achieve stronger security and more importantly, the leaked information of the encrypted vault eventually makes it less secure.*

**The security of our scheme.** We only show the password-similarity attack and the hybrid attack with Feature Overlap in Fig. 4 and Table 4, since the feature performs the best for attacks among the four features demonstrated in Section 4.3. The experimental results for other features are given in Appendix E.

As shown in Fig. 4 and Table 4, the hybrid attack delivers the best performance, where the average rank  $\bar{r}$  and the accuracy  $\alpha$  are 42% and 58%, respectively. Compared to the existing schemes with 85%–94% accuracy and 6%–15% average rank, our scheme brings 2.8–7.5 times cost of online verifications to attackers. Since online verifications can be quickly detected and prevented [16, 20, 40], *our scheme does make a significant improvement on resisting distinguishing attacks in the single-leakage case.*

The decreased cracked proportions also illustrate the security improvement. As shown in Table 6, the existing schemes

suffer from 33%–55% (i.e.,  $F(0)$ ) real vaults cracking via one guess, this value is only 2% for our scheme, which decreases the harm by 93%–96%.

## 5.4 Further Discussion

**Other experiments.** We also evaluate the security of honey vault schemes against the password policy attack and the intersection attack. The experimental results are trivial: the attacks completely breaks the existing schemes, but are resisted by ours (see Appendix C for the intersection attack and the full version of this paper for the password policy attack).

**Limitation on the dataset.** The vault dataset, Pastebin, we used, is not leaked from real vault applications and its size is relatively small (see Section 5.2). Although it is well-studied and used in [12, 13, 17] for security evaluation, the quality of it may yield some bias in our experimental results. Nevertheless, our experiments still demonstrate the insecurity of the existing honey vault schemes [12, 13, 17]. Furthermore, the quality of the dataset does not affect some important conclusions: 1) our construction with the incremental update mechanism can resist intersection attacks; 2) an adaptive scheme leaks extra information of the encrypted vault and is less secure than its static variant.

In general, a dataset with better data quality may help an attacker to more precisely model the real vault/password distribution and more effectively distinguish real and decoy vaults against honey vault schemes (including ours). On the other hand, such a dataset may benefit the design of vault models. Via our generic construction roadmap, using a more accurate multi-similar-password model can generate more plausible-looking decoys and the update security will be maintained. We note that designing a model and cracking it is not a cat-and-mouse game. If one can design a vault model that precisely captures most of the vaults in the vault space, then arbitrary attackers, even with better-quality datasets, will have little advantage in distinguishing real and decoy vaults.

**More powerful attacks.** There may be some other information that could be used (as pre-knowledge) to launch attacks. Personal information may be one of the options, since users may create passwords based on name, birthday, phone number, email and username [31, 48]. Attackers may identify a real vault with a higher probability. e.g., by checking if the passwords in the vault match personal information. To resist this type of attack, we may consider using a conditional probability model (e.g., the Personal-PCFG model [31]) that is able to characterize the real vault distribution under the condition of the provided information.

The size of a vault may also provide an extra advantage for attackers in launching online verification. With more accounts on different websites, attackers may be allowed to launch more online verifications on these websites. We do not consider this advantage and leave it as future work.

## 6 Conclusion and Future Work

We propose a generic construction and further an incremental update mechanism for honey vault schemes. The update mechanism enables the vault scheme to achieve the update security, i.e., *resisting intersection attacks* in the multi-leakage case. We instantiate our scheme with a well-designed multi-similar-password model. Our evaluation with real-world datasets shows that compared with the existing schemes, our instance achieves a significant improvement on security *against (traditional) distinguishing attacks* in the single-leakage case.

Our work may also benefit other research topics. For example, the incremental update mechanism can be used for other HE applications, the multi-similar-password model may benefit password guessing attacks and password strength meters. We leave these as future work.

## Acknowledgment

The authors are grateful to the anonymous reviewers and the shepherd, David Freeman, for their invaluable comments that highly improve the completeness of the paper. We also give our special thanks to Qianchen Gu, Zhixiong Zheng, Jiahong Yang, Xiaoxi He and Jiahong Xie for their insightful suggestions and invaluable help. This research is supported by National Key R&D Program of China (2020YFB1805400), National Natural Science Foundation of China (62072010), and European Union’s Horizon 2020 research and innovation programme under grant agreement No. 952697 (ASSURED).

## References

- [1] Have i been pwned? <https://haveibeenpwned.com>.
- [2] How do I view username, password, and note history for sites? <https://support.logmeininc.com/lastpass/help/how-do-i-view-username-password-and-note-history-for-sites>.

- [3] LastPass technical whitepaper. <https://enterprise.lastpass.com/wp-content/uploads/LastPass-Technical-Whitepaper-3.pdf>.
- [4] 1Password security design, January 2019. <https://1password.com/files/1Password-White-Paper.pdf>.
- [5] Issue 1930: lastpass: bypassing do\_popupregister() leaks credentials from previous site, October 2019. <https://bugs.chromium.org/p/project-zero/issues/detail?id=1930>.
- [6] 1Password. 1password backups. <https://support.1password.com/backups/>.
- [7] Jeremiah Blocki, Ben Harsha, and Samson Zhou. On the economics of offline password cracking. In *IEEE S&P 2018*, pages 35–53.
- [8] Hristo Bojinov, Elie Bursztein, Xavier Boyen, and Dan Boneh. Kamouflage: Loss-resistant password management. In *ESORICS 2010*, pages 286–302. Springer.
- [9] Dan Boneh, Henry Corrigan-Gibbs, and Stuart Schechter. Balloon hashing: A memory-hard function providing provable protection against sequential attacks. In *ASIACRYPT 2016*, pages 220–248. Springer.
- [10] Joseph Bonneau and Stuart Schechter. Towards reliable storage of 56-bit secrets in human memory. In *USENIX Security 2014*, pages 607–623.
- [11] Rahul Chatterjee, Anish Athayle, Devdatta Akhawe, Ari Juels, and Thomas Ristenpart. password typos and how to correct them securely. In *IEEE S&P 2016*.
- [12] Rahul Chatterjee, Joseph Bonneau, Ari Juels, and Thomas Ristenpart. Cracking-resistant password vaults using natural language encoders. In *IEEE S&P 2015*, pages 481–498.
- [13] Haibo Cheng, Zhixiong Zheng, Wenting Li, Ping Wang, and Chao-Hsien Chu. Probability model transforming encoders against encoding attacks. In *USENIX Security 2019*, pages 1573–1590.
- [14] Anupam Das, Joseph Bonneau, Matthew Caesar, Nikita Borisov, and XiaoFeng Wang. The tangled web of password reuse. In *NDSS 2014*, pages 1–15.
- [15] Dropbox. File version history overview. <https://help.dropbox.com/files-folders/restore-delete/version-history-overview>.
- [16] David Freeman, Sakshi Jain, Markus Dürmuth, Battista Biggio, and Giorgio Giacinto. Who are you? a statistical approach to measuring user authenticity. In *NDSS 2016*, pages 1–15.
- [17] Maximilian Golla, Benedict Beuscher, and Markus Dürmuth. On the security of cracking-resistant password vaults. In *ACM CCS 2016*, pages 1230–1241.
- [18] Google. Protect your accounts from data breaches with password checkup. <https://security.googleblog.com/2019/02/protect-your-accounts-from-data.html>.
- [19] Amber Gott. LastPass security notification, March 2017. <https://blog.lastpass.com/2017/03/important-security-updates-for-our-users.html/>.
- [20] Paul A Grassi, James L Fenton, Elaine M Newton, Ray A Perlner, Andrew R Regenscheid, William E Burr, and Justin P

- Richer. Nist special publication 800-63b. Digital identity guidelines: Authentication and lifecycle management. *Bericht, NIST*, 2017.
- [21] Joseph Jaeger, Thomas Ristenpart, and Qiang Tang. Honey encryption beyond message recovery security. In *EUROCRYPT 2016*, pages 758–788, 2016.
- [22] Ari Juels and Thomas Ristenpart. Honey encryption: Security beyond the brute-force bound. In *EUROCRYPT 2014*, pages 293–310. Springer.
- [23] Burt Kaliski. PKCS #5: Password-based cryptography specification version 2.0. 2000.
- [24] Mathias Karlsson. How I made LastPass give me all your passwords, July 2016. <https://labs.detectify.com/2016/07/27/how-i-made-lastpass-give-me-all-your-passwords/>.
- [25] Jason Kincaid. Dropbox security bug made passwords optional for four hours, June 2011. <https://techcrunch.com/2011/06/20/dropbox-security-bug-made-passwords-optional-for-four-hours/>.
- [26] Eugene F Krause. *Taxicab geometry: An adventure in non-Euclidean geometry*. Courier Corporation, 1986.
- [27] Russell WF Lai, Christoph Egger, Manuel Reinert, Sherman SM Chow, Matteo Maffei, and Dominique Schröder. Simple password-hardened encryption services. In *USENIX Security 2018*, pages 1405–1421.
- [28] Michael Levandowsky and David Winter. Distance between sets. *Nature*, 234(5323):34, 1971.
- [29] Vladimir I Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, pages 707–710, 1966.
- [30] Dave Lewis. iCloud data breach: Hacking and celebrity photos, September 2014. <https://www.forbes.com/sites/davelewis/2014/09/02/icloud-data-breach-hacking-and-nude-celebrity-photos/>.
- [31] Yue Li, Haining Wang, and Kun Sun. A study of personal information in human-chosen passwords and its security implications. In *IEEE INFOCOM 2016*, pages 1–9.
- [32] Zhiwei Li, Warren He, Devdatta Akhawe, and Dawn Song. The emperor’s new password manager: Security analysis of web-based password managers. In *USENIX Security 2014*, pages 465–479.
- [33] Sanam Ghorbani Lyastani, Michael Schilling, Sascha Fahl, Sven Bugiel, and Michael Backes. Better managed than memorized? studying the impact of managers on password strength and reuse. In *USENIX Security 2018*, pages 203–220.
- [34] Jerry Ma, Weining Yang, Min Luo, and Ninghui Li. A study of probabilistic password models. In *IEEE S&P 2014*, pages 538–552.
- [35] William Melicher, Blase Ur, Sean M Segreti, Saranga Komanduri, Lujo Bauer, Nicolas Christin, and Lorrie Faith Cranor. Fast, lean, and accurate: Modeling password guessability using neural networks. In *USENIX Security 2016*, pages 175–191.
- [36] Bijeeta Pal, Tal Daniel, Rahul Chatterjee, and Thomas Ristenpart. Beyond credential stuffing: Password similarity models using neural networks. In *IEEE S&P 2019*, pages 814–831.
- [37] Dario Pasquini, Ankit Gangwal, Giuseppe Ateniese, Massimo Bernaschi, and Mauro Conti. Improving password guessing via representation learning. In *IEEE S&P 2021*, pages 265–282.
- [38] Sarah Pearman, Shikun Aerin Zhang, Lujo Bauer, Nicolas Christin, and Lorrie Faith Cranor. Why people (don’t) use password managers effectively. In *SOUPS 2019*, pages 319–338.
- [39] Colin Percival. Stronger key derivation via sequential memory-hard functions. *Self-published*, pages 1–16, 2009.
- [40] Benny Pinkas and Tomas Sander. Securing passwords against dictionary attacks. In *ACM CCS 2002*, pages 161–170.
- [41] Niels Provos and David Mazieres. A future-adaptable password scheme. In *USENIX ATC 1999*, pages 81–91.
- [42] Joe Siegrist. LastPass security notification, May 2011. <https://blog.lastpass.com/2011/05/lastpass-security-notification.html/>.
- [43] Joe Siegrist. LastPass hacked – identified early & resolved, July 2015. <https://blog.lastpass.com/2015/06/lastpass-security-notice.html/>.
- [44] David Silver, Suman Jana, Dan Boneh, Eric Yawei Chen, and Collin Jackson. Password managers: Attacks and defenses. In *USENIX Security 2014*, pages 449–464.
- [45] Karen Turner. Hacked dropbox login data of 68 million users is now for sale on the dark web, September 2016. <https://www.washingtonpost.com/news/the-switch/wp/2016/09/07/hacked-dropbox-data-of-68-million-users-is-now-for-sale-on-the-dark-web/>.
- [46] Rafael Veras, Christopher Collins, and Julie Thorpe. On the semantic patterns of passwords and their security impact. In *NDSS 2014*, pages 1–16.
- [47] Ding Wang, Haibo Cheng, Ping Wang, Jeff Yan, and Xinyi Huang. A security analysis of honeywords. In *NDSS 2018*, pages 1–15.
- [48] Ding Wang, Zijian Zhang, Ping Wang, Jeff Yan, and Xinyi Huang. Targeted online password guessing: An underestimated threat. In *ACM CCS 2016*, pages 1242–1254.
- [49] Matt Weir, Sudhir Aggarwal, Breno de Medeiros, and Bill Glodek. Password cracking using probabilistic context-free grammars. In *IEEE S&P 2009*, pages 391–405.
- [50] Jeff Yan, Alan Blackwell, Ross Anderson, and Alasdair Grant. Password memorability and security: Empirical results. *IEEE Secur. Priv.*, 2(5):25–31, 2004.

## A Our Single-Similar-Password Model and Its Conditional Encoder

**Our design.** For simplicity, we only consider the most popular reuse habit, i.e., head or tail modifications. Specifically, in our model, the new password  $pw_{\text{new}}$  can be generated from an old one  $pw_{\text{old}}$  in the following ways:

1. Direct reuse, i.e.,  $pw_{\text{new}} = pw_{\text{old}}$ .
2. Modify the characters in the tail of  $pw_{\text{old}}$ . The modification operations include deleting, adding and deleting-then-adding. For example, given  $pw_{\text{old}} = \text{“password!”}$ ,  $pw_{\text{new}}$  may be “password” (deleting the last charac-



ter “!”), “password!\*” (adding a character “\*”), or “password\*@” (deleting “!” and then adding “\*@”).

3. Modify the characters in the head of  $pw_{old}$  with the same modification operations as above.
4. Modify the head characters and then the tail characters.

In addition, we find that users prefer to directly reuse passwords (i.e., without any modifications) if they already have many passwords. This is intuitive, due to the memory limitation of humans. Therefore, the probability of direct reuse increases as the number of old passwords increases. However, current models [36, 48] only use a static probability for direct reuse, and therefore are not suitable for our construction of the multi-similar-password model. To address this issue, we set an adaptive probability for direct reuse and quantify the probability based on the real-world vault dataset, Pastebin. We find that the proportion of the same password pairs (i.e., direct reuse) in the similar password pairs varies little with vault size. Assuming the proportion is a constant (denoted as  $\alpha$ ) independent of vault size (i.e., for different vault sizes, an arbitrary similar pair is a same pair with a constant probability  $\alpha$ ), we have that the probability of direct reuse is

$$\frac{i \times \alpha}{i \times \alpha + 1 - \alpha},$$

where  $i$  is the number of previous passwords (i.e.  $pw_{new}$  is the  $i + 1$ -th password). When training, we use the average proportion in the training set for  $\alpha$ . Except this probability of direct reuse, other probabilities (e.g. those of modifying characters) are simply counted from the training dataset.

To show the details of character modification in our model, we give an example with “password!” and “password@1” as the old and new passwords. Clearly, “password@1” is generated from “password!” by deleting “!” and adding “@1”. Then  $\Pr_{SSPM}(\text{password}@1 | \text{password}!) = \Pr_{DR}(\text{False}) \times \Pr_M(\text{Tail}) \Pr_T(\text{Deleting-then-adding}) \Pr_{TDN}(1) \Pr_{TAN}(2) \times \Pr_{TAC}(@) \Pr_{TAC}(1)$ . Here,  $\Pr_{DR}(\text{False})$ ,  $\Pr_M(\text{Tail})$ , and  $\Pr_T(\text{Deleting-then-adding})$  are the probabilities of not direct reuse, modifying tail, and deleting-then-adding tail characters, respectively;  $\Pr_{TDN}(1)$  and  $\Pr_{TAN}(2)$  are the probabilities of deleting 1 tail character and adding 2 tail characters, respectively;  $\Pr_{TAC}(@)$  and  $\Pr_{TAC}(1)$  are the probabilities of adding characters “@” and “1” to the tail, respectively.

Several details should be carefully considered in our model. Let  $l_{old}$  be the length of the old password  $pw_{old}$ ,  $l_{HD}$ ,  $l_{TD}$ ,  $l_{HA}$ , and  $l_{TA}$  be the character numbers of head deleting, tail deleting, head adding, and tail adding, respectively. Then, the length of the new password  $pw_{new}$  is  $l_{new} = l_{old} - l_{HD} - l_{TD} + l_{HA} + l_{TA}$  and the longest common substring length  $l_{LCSStr}$  of  $pw_{old}$  and  $pw_{new}$  is  $l_{old} - l_{HD} - l_{TD}$ . Because  $\frac{1}{2}l_{old} \leq l_{LCSStr} \leq 2l_{old}$ , it holds that  $l_{HD} \leq \frac{1}{2}l_{old}$ ,  $l_{TD} \leq \frac{1}{2}l_{old} - l_{HD}$ ,  $l_{HA} \leq 2(l_{old} - l_{HD} - l_{TD})$ , and  $l_{TA} \leq 2(l_{old} - l_{HD} - l_{TD}) - l_{HA}$ . Therefore, when calculating  $\Pr_{HDN}$ ,  $\Pr_{TDN}$ ,  $\Pr_{HAN}$  and  $\Pr_{TAN}$ , all invalid values in the tables should be excluded and meanwhile, the probabilities of the remaining values should be

normalized. Note this process is the same as the pruning method [13]. As a result, a decoy seed can be always decoded to a valid vault. Furthermore, if at least one character in the head (or tail) is deleted, then the first head-added (or tail-added) character cannot be the same character as the old one (but other added characters can be identical). This helps us reduce the ambiguity of our model. Similar changes should be applied to  $\Pr_{HAC}$  and  $\Pr_{TAC}$ .

With the above designs, our model significantly reduces ambiguity but still cannot eliminate it. This is due to the non-uniqueness of the longest common substring (note the same longest common substrings on different positions are treated as two different ones). For instance, the password “aaaaa” can be modified to “aaaa” in two different ways: deleting the first or the last character. In this case, the probability  $\Pr_{SSPM}(pw_{old} | pw_{new})$  is defined as the total probability of all modifying methods.

**Conditional encoder for our model.** We use the method proposed in Section 3.2 to convert this model to a conditional encoder. As discussed above, the conditional encoder needs to parse all the longest common substrings of two passwords. With a generalized suffix tree, this operation can be done in  $O(l_1 l_2)$  time, where  $l_1, l_2$  are the lengths of the two passwords, respectively. Other operations of the encoder are simple and fast. Thus, our encoder is efficient for real applications.

If one sets our model to characterize more transformation rules of password reuse, then the resulting encoder will suffer from time complexity. This is the reason why we prefer to keep our model simple.

## B Conditional Encoder for Our Multi-Similar-Password Model

In our multi-similar-password model, the new password  $pw_{i+1}$  is generated by reusing  $pw_{i'}$  ( $1 \leq i' \leq i$ ) or a brand new selection. This means a valid generating path has the form  $(g, r_1, r_2, \dots)$ , where  $g$  represents the generating path:  $i'$  for reusing  $pw_{i'}$  and 0 for otherwise. Note that: if  $g = 0$ , then  $(r_k)_k$  must be a generating path in the single-password model; if  $g = i'$ , then  $(r_k)_k$  must be a generating path in the single-similar-password model under the condition of  $pw_{i'}$ . Using the model-to-encoder transformation proposed in Section 3.2, we can construct the following conditional encoder for our multi-similar-password model.

To encode  $pw_{i+1}$  with  $i$  given old passwords  $(pw_{i'})_{i'=1}^i$ , the encoder works as

1. Calculate  $\frac{1-f(i)}{i} \Pr_{SSPM}(pw_{i+1} | pw_{i'})$  for  $1 \leq i' \leq i$  and  $f(i) \Pr_{SPM}(pw_{i+1})$ .
2. According to the above probabilities, choose a generating path of  $pw_{i+1}$ . If the generating path  $s$  is modifying  $pw_{i'}$ , set  $g = i'$ , otherwise, set  $g = 0$ .
3. Encode  $g$  by the IS-DTE for the distribution of  $g$  (the probability is  $f(i)$  for 0 and  $\frac{1-f(i)}{i}$  for 1 to  $i$ ).

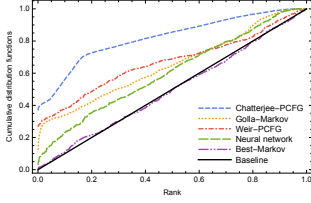


Figure 5: RCDFs for different models trained with RockYou

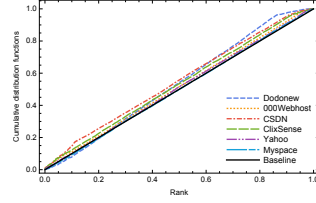


Figure 6: RCDFs for Best-Markov trained with different datasets

Table 7: The average rank  $\bar{r}$  for Markov model under the single-password attack

Normalization		End-symbol			Distribution-based		
Order		3	4	5	3	4	5
Pseudocount	$10^0$	31%	35%	39%	33%	37%	39%
	$10^{-1}$	30%	37%	43%	34%	38%	44%
	$10^{-2}$	31%	36%	44%	33%	40%	48%
	$10^{-3}$	32%	35%	42%	33%	39%	<b>50%</b>
	$10^{-4}$	30%	39%	45%	34%	42%	48%
	$10^{-5}$	32%	36%	46%	32%	42%	47%

- If  $g = 0$ , encode  $pw_{i+1}$  (the remaining rules in generating path  $s$ ) by the encoder (done by Cheng et al.’s transformation [13]) for our single-password model, otherwise, encode  $pw_{i+1}$  by the conditional encoder (made by our extended transformation) for our single-similar-password model with the given old password  $pw_g$ .
- Concatenate these seeds, pad the concatenation to a fixed-length seed with random bits, and output the seed.

## C Intersection Attacks Against Honey Vault Schemes

We evaluate the security of honey vault schemes against intersection attacks with the real-world datasets.

**Experimental settings.** For each vault (with a size larger than 2) in the vault dataset Pastebin, we randomly shuffle the passwords in the vault and treat the last password as a new added one. In this way, we get the old and new versions for each (real) vault (here, the old version is the vault without the last password). Then we use the honey vault schemes to generate the decoys for these real vaults (note the decoys have respective two versions as well). The rest of the experiment settings are the same as those in Section 5.

**Intersection attacks.** We leverage a trivial intersection attack for the evaluation. This attack only leverages the similarity between the old and new versions of vaults, but not considers the difference between the real vault distribution and the vault model. For each candidate vault  $V_i$  with its old and new versions  $V_i^o, V_i^n$ , the priority function  $p_{TIA}(V_i)$  of the trivial intersection attack is equal to 1 if  $V_i^o$  is the same as  $V_i^n$  except for the last password, otherwise, 0. In other words, the

Table 8: The average rank  $\bar{r}$  for single-password models under the single-password attack

Model	Chatterjee-PCFG	Golla-Markov	Weir-PCFG	Neural network	Best-Markov <sup>1</sup>
$\bar{r}$	18%	35%	33%	40%	50%

<sup>1</sup> Best-Markov is the 5-order Markov model using distribution-based normalization and Laplace smoothing with the pseudocount of 0.001.

Table 9: The average rank  $\bar{r}$  for Best-Markov trained with different datasets under the single-password attack

Dataset	Dodonew	000Webhost	CSDN	ClixSense	Yahoo	Myspace
$\bar{r}$	47%	48%	45%	47%	49%	50%

attack directly excludes these vaults of which two versions are not similar.

**Experimental results.** The intersection attack can directly tell the real vault with 100% accuracy for all existing honey vault schemes. This is because the two versions of each decoy vault are randomly generated and there is only a very small probability that the two versions are similar. In contrast, our scheme generates the new version of each decoy vault by adding a new password to the old version. Therefore, it can resist the intersection attack.

## D Evaluating Single-Password Models

To instantiate our construction with a good single-password model, we evaluate the existing models with the single-password attack.

**Existing single-password models.** We here evaluate Chatterjee-PCFG [12], Golla-Markov [17], Weir-PCFG [49], the neural network model [35], and Markov models [34] using different methods. The neural network model proposed by Melicher et al. [35] is the same as 10-order Markov model except that the transition probabilities are calculated by recurrent neural networks. The performance of the model heavily depends on its parameter settings. Here we use the default settings in the example configuration on GitHub. Similar to the neural network model, the performance of Markov models relies on its normalization and smoothing methods. We note that in the Markov models with distribution-based normalization, we train independent Markov models for passwords with different lengths, which is not considered in [34] but [17]. We employ this operation because it can make a significant improvement on the performance of Markov models. So, we only present the experiments with this operation.

**Experimental settings.** We randomly choose 50% passwords from RockYou as the training set for PMTEs and attacks, while randomly selecting  $10^4$  passwords from the remaining part of RockYou as real passwords. Apart from that, other experimental settings are the same as those for honey vaults.

**Experimental results.** As shown in Fig. 5, Tables 7 and

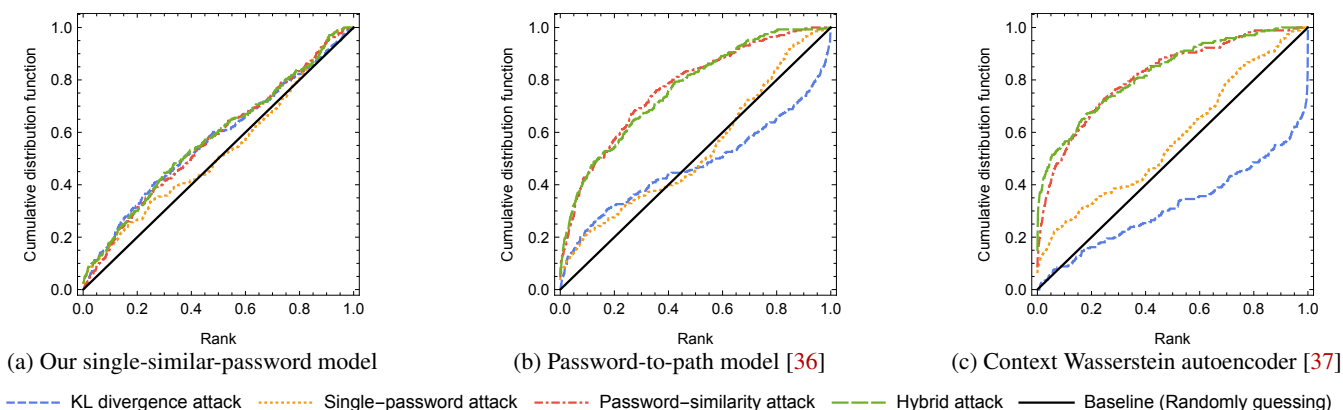


Figure 7: RCDFs for our honey vault scheme using different single-similar-password models under attacks (with the best performance similar meter)

Table 10: The average rank  $\bar{r}$  for our honey vault scheme with different single-similar-password models

Attack	Ours	Password-to-path model [36]	Context Wasserstein autoencoder [37]
KL divergence	42%	52%	68%
Single password	48%	47%	42%
Password similarity (Levenshtein)	45%	24%	23%
Password similarity (LCS)	47%	31%	<b>18%</b>
Password similarity (Manhattan)	46%	<b>23%</b>	24%
Password similarity (Overlap)	<b>43%</b>	37%	22%
Hybrid (Levenshtein)	44%	26%	22%
Hybrid (LCS)	46%	28%	19%
Hybrid (Manhattan)	43%	<b>24%</b>	22%
Hybrid (Overlap)	<b>42%</b>	35%	<b>17%</b>

8, the 5-order Markov models using distribution-based normalization and Laplace smooth with pseudocount of  $10^{-3}$  can guarantee the expected security, where the average rank is 50%. We denote these models as *Best-Markov*. We also achieve good performance for Best-Markov via using other password datasets including Dodonew, 000Webhost, CSDN, ClixSense, Yahoo, and Myspace. Note these datasets are extensively used in password researches, e.g., [34, 47, 49]. As shown in Fig. 6 and Table 9, the average ranks of Best-Markov all approach to the expected value ( $\approx 50\%$ ) and meanwhile, the RCDFs are all close to the baseline. Due to the good performance of Best-Markov, we will use it in our vault model.

## E Evaluating Single-Similar-Password Models

To choose a single-similar-password model for our multi-similar-password model, we evaluate two existing models, the password-to-path model (pass2path) [36] and the context Wasserstein autoencoder (CWAE) [37] along with our simple model (Appendix A) under our proposed attacks. Wang et al. [48] also propose a design, but we do not consider it be-

cause: 1) pass2path always outperforms this model on target password guessing [36]; 2) Wang et al. do not open-source the code, which brings difficulty in comparison.

**Experimental settings.** The settings are the same as those in Section 5, except the following.

1. *Training.* We leverage pass2path and CWAE which are trained by the authors and provided on GitHub. We do not train these two models with the real dataset Pastebin, because the dataset is too small for them and cannot provide sufficient training. We note that our model is simple and can be trained with a small-scale dataset.
2. *Password-similarity features.* Pass2path is training from the password pairs  $(pw_1, pw_2)$  satisfying that the Levenshtein distance between  $pw_1$  and  $pw_2$  is not more than 5. So we use this feature as Feature M in the password-similarity attack. As for CWAE, we still use Feature LCSStr as Feature M.
3. *Encoder.* We directly sample passwords from pass2path and CWAE to generate decoy vaults in the experiments, instead of implementing their encoders. This is because these models yield heavy time complexity in encoding. For example, in pass2path, “Password” can be generated from “password” 1) by capitalizing “p” or 2) by deleting “p” then inserting “P” or 3) via other paths. To resist encoding attacks, the encoder needs to parse all generating paths [13], yielding exponential time complexity.

**Experimental results.** As shown in Fig. 7 and Table 10, our simple design performs the best: the average rank  $\bar{r}$  is not less than 42% under any attacks. This means no attacks achieve more than 58% accuracy in distinguishing real and decoy vaults. Both pass2path and CWAE cannot generate plausible-looking decoys: targeting pass2path, the password similarity attack with the Manhattan meter achieves 77% accuracy; targeting CWAE, the hybrid attack with the Overlap meter achieves 83% accuracy. Since our simple model performs well on decoy generating, it may also perform well on password guessing. We leave this as future work.