# Spaced Repetition and Mnemonics Enable Recall of Multiple Strong Passwords

Jeremiah Blocki
Computer Science Department
Carnegie Mellon University
Pittsburgh, PA 15213
Email: jblocki@cs.cmu.edu

Saranga Komanduri
Computer Science Department
Carnegie Mellon University
Pittsburgh, PA 15213
Email: sarangak@cs.cmu.edu

Lorrie Cranor
Computer Science Department
Carnegie Mellon University
Pittsburgh, PA 15213
Email: lorrie@cs.cmu.edu

Anupam Datta
Computer Science Department
Carnegie Mellon University
Pittsburgh, PA 15213
Email: danupam@cs.cmu.edu

arXiv:1410.1490v1 [cs.CR] 6 Oct 2014

*Abstract*—We report on a user study that provides evidence that spaced repetition and a specific mnemonic technique enable users to successfully recall multiple strong passwords over time. Remote research participants were asked to memorize 4 Person-Action-Object (PAO) stories where they chose a famous person from a drop-down list and were given machine-generated random action-object pairs. Users were also shown a photo of a scene and asked to imagine the PAO story taking place in the scene (e.g., Bill Gates—swallowing—bike on a beach). Subsequently, they were asked to recall the action-object pairs when prompted with the associated scene-person pairs following a spaced repetition schedule over a period of $100+$ days. While we evaluated several spaced repetition schedules, the best results were obtained when users initially returned after 12 hours and then in $1.5\times$ increasing intervals: $77.1\%$ of the participants successfully recalled all $4$ stories in $9$ tests over a period of $102$ days. Much of the forgetting happened in the first test period ($12$ hours): on average $94.9\%$ of the participants who had remembered the stories in earlier rounds successfully remembered them in subsequent rounds. These findings, coupled with recent results on naturally rehearsing password schemes, suggest that $4$ PAO stories could be used to create usable and strong passwords for $14$ sensitive accounts following this spaced repetition schedule, possibly with a few extra upfront rehearsals. In addition, we find statistically significant evidence that initially ($8$ tests over $64$ days) users who were asked to memorize $4$ PAO stories outperform users who are given $4$ random action-object pairs, but eventually ($9$ tests over $128$ days) the advantage is not significant. Furthermore, there is an interference effect across multiple PAO stories: the recall rate of $100\%$ for participants who were asked to memorize $1$ or $2$ PAO stories is significantly better than that for $4$ PAO stories. These findings yield concrete advice for improving constructions of password management schemes and future user studies.

## I. INTRODUCTION

Passwords are currently the dominant form of human authentication over the Internet despite many attempts to replace them [1]. A typical internet user has the complex task of creating and remembering passwords for many different accounts. Users struggle with this task, adopting insecure password practices [2]–[5] or frequently having to reset their passwords. Yet research on human memory provides reason for optimism. Specifically, *spaced repetition*—a memorization technique that incorporates increasing intervals of time between subsequent

review of previously learned material—has been shown to be effective in enabling recall in a wide variety of domains [6]–[10]. Similarly, *mnemonic* techniques that provide multiple semantic encodings of information (e.g., as stories and images) also significantly help humans recall information [10], [11].

We report on a user study that provides evidence that spaced repetition and mnemonics enable users to successfully recall *multiple* strong passwords over time. The study is inspired by a recent result on *naturally rehearshing password schemes* [12] that rely on spaced repetition and a specific Person-Action-Object (PAO) mnemonic technique to design a scheme to create and maintain multiple strong passwords. As a core component of the study, remote research participants were asked to memorize 4 Person-Action-Object (PAO) stories where they chose a famous person from a drop-down list and were given machine-generated random action-object pairs. Users were also shown a photo of a scene and asked to imagine the PAO story taking place in the scene (e.g., Bill Gates—swallowing—bike on a beach). Subsequently, they were asked to recall the action-object pairs (e.g., swallowing—bike) when prompted with the associated scene-person pairs (e.g., Bill Gates—beach) following a spaced repetition schedule over a period of $100+$ days. We designed the study to seek answers to the following questions:

- Do users who follow spaced repetition schedules successfully recall multiple PAO stories and, if so, which schedules work best?
- Does the PAO mnemonic technique improve recall over random action-object pairs?
- Is there an interference effect when users are asked to memorize multiple PAO stories?

We summarize our key findings and discuss their implications for password management below. First, while we evaluated several spaced repetition schedules, the best results were obtained under the schedule in which users initially returned after 12 hours and then in $1.5\times$ increasing intervals: $77.1\%$ of the participants successfully recalled all $4$ stories in $9$ tests over a period of $102$ days. Much of the forgetting happened in the first test period (the first 12 hours): on average $94.9\%$ of the participants who had remembered the stories in earlier rounds successfully remembered them in subsequent

rounds. These findings, coupled with the results of Blocki et al. [12], suggest that 4 PAO stories could be used to create and maintain usable and strong passwords for up to 14 accounts following this spaced repetition schedule, possibly with a few extra upfront rehearsals. The finding that much of the forgetting happens in the first test period robustly held in all the spaced repetition schedules that we experimented with. Another implication of this finding is that password expiration policies [13] negatively impact usability by forcing users to return to the highest rehearsal effort region of memorizing a password. Furthermore, they are unnecessary for strong passwords (see Section II).

Second, we find statistically significant evidence that initially (8 tests over 64 days) users who were asked to memorize 4 PAO stories outperform users who are given 4 random action-object pairs, but eventually (9 tests over 128 days) the advantage is not significant. This finding is consistent with the previous finding in that much of the forgetting happens in the early rounds and in those rounds the PAO mnemonic technique helps significantly with recall.

Third, we find a statistically significant interference effect across multiple PAO stories. Specifically, the recall rate of 100% for participants who were asked to memorize 1 or 2 PAO stories is significantly better than the rate for participants who were asked to memorize 4 PAO stories. The interference effect is strong: it continues to be statistically significant even if we only count a participant with 4 PAO stories as failing if they forgot their first (or first two) action-object pair(s). This finding has several implications for password management. Further studies are needed to discover whether the interference effect is alleviated if users memorize multiple PAO stories following a staggered schedule in which they memorize 2 stories at a time. To accomodate this user model, we also need new constructions for naturally rehearsing password schemes in which passwords can be constructed even when not all PAO stories are memorized upfront (see Section VI for a concrete open problem). At the same time, the perfect recall rate for 1 or 2 PAO stories suggests that they could serve as a mechanism for strengthening existing passwords over time. This conclusion is similar to the conclusion of a related study of Bonneau and Schechter [14] (although there are significant differences between the two studies that we discuss in Section V).

*Organization.* Section II briefly reviews the password management scheme of Blocki et al. [12], and the security of the associated passwords consisting of random action-object pairs. Section III presents the design of our user study. Section IV describes the results of the study. Section V describes related work. Finally, Section VI concludes with a discussion of the implications of these results for password management and suggestions for future work.

## II. BACKGROUND: SHARED CUES

In this section we analyze the security of passwords consisting of random action-object pairs (Section II-A) and overview the Shared Cues password management scheme of Blocki et

al. [12]. In Section II-C we consider a variation of the Shared Cues password management scheme which only requires the user to memorize four PAO stories to form 14 strong passwords.

### A. Security Against Offline Attacks

Any adversary who has obtained the cryptographic hash of a user's password can mount an automated brute-force attack to crack the password by comparing the cryptographic hash of the user's password with the cryptographic hashes of likely password guesses. This attack is called an offline dictionary attack, and there are many password crackers that an adversary could use [15]. Offline dictionary attacks against passwords are powerful and commonplace [16]. Adversaries have been able to compromise servers at large companies (e.g., Zappos, LinkedIn, Sony, Gawker [17]–[22]) resulting in the release of millions of cryptographic password hashes[1].

In our study each action is chosen uniformly at random from a list of 92 actions and each object is chosen from a list of 96 objects. Thus a randomly chosen action-object pair has $\approx 13.11$ bits of entropy (approximately equivalent to the 13.29 bits of entropy in a randomly chosen 4-digit pin number). We assume that passwords are hashed using a cryptographic password hash function $\mathbb{H}$ and we let $\mathbb{C}(\mathbb{H})$ denote the cost of evaluating the cryptographic hash function one time. Bonneau and Schechter used data on the Bitcoin mining economy to estimate that $\mathbb{C}(\mathbb{H}) \approx \$2^{-50.07}$ for the SHA-256 hash function, and they estimate that with iterated password hashing we can increase this cost to $\mathbb{C}(\mathbb{H}) \approx \$2^{-26.07}$ if we are willing to wait approximately two seconds to compute $\mathbb{H}$ during authentication[2]. We note that on a computer with $2^d$ cores (e.g., GPU) the value of $\mathbb{C}(H)$ can be adjusted by a factor of $2^d$ without increasing the authentication time as follows: select a random salt value $s \in \{0,1\}^*$ as well as a random value $x \sim \{0,1\}^d$, compute $\mathbb{H}(pw, x, s)$ and store $(s, \mathbb{H}(pw, x, s))$. Thus, if authentication were performed on a GPU with $1,024$ cores and we were willing to wait approximately two seconds for authentication we could increase $\mathbb{C}(H) \approx \$2^{-16.07} \approx \$1.46 \times 10^{-5}$. If a password is chosen uniformly at random from a space of size $N$ then the adversary's expected cost of an offline attack is $N\mathbb{C}(\mathbb{H})/2$. Table I shows the expected cost of an offline attack. A password consisting of two secret action-object pairs would be sufficiently strong to protect many accounts as long as the value $\mathbb{C}(\mathbb{H}) \times 10^6$ is least \$1 — Symantec reported that compromised passwords are sold for between \$4 and \$30 on the black market [24], and a password consisting of three action-object pairs would be sufficiently strong to protect most high value accounts.

### B. Shared Cues Password Management Scheme

Our user study is partially motivated by the Shared Cues password management scheme of Blocki et al. [12]. In their scheme the user memorizes random PAO stories, and forms

---

[1]In a few of these cases [21], [22] the passwords were stored in the clear.
[2]Cryptographic password hash functions like SCRYPT or BCRYPT [23] use similar ideas to increase $\mathbb{C}(\mathbb{H})$ .

| | # Action-Object Pairs in Password | | | |
|---|---|---|---|---|
| $\mathbb{C}(\mathbf{H})$ | One | Two | Three | Four |
| $\$10^{-5}$ | $\$4.4 \times 10^{-2}$ | $\$390$ | $\$3.4 \times 10^{6}$ | $\$3.0 \times 10^{10}$ |
| $\$10^{-6}$ | $\$4.4 \times 10^{-3}$ | $\$39$ | $\$3.4 \times 10^{5}$ | $\$3.0 \times 10^{9}$ |
| $\$10^{-7}$ | $\$4.4 \times 10^{-4}$ | $\$3.9$ | $\$3.4 \times 10^{4}$ | $\$3.0 \times 10^{8}$ |

TABLE I: Expected Cost of an Offline Attack

his passwords by appending the secret action(s) and object(s) from different stories together.

*Person-Action-Object Stories.* A user who adopts the Shared Cues password management scheme [12] first memorizes several randomly generated Person-Action-Object (PAO) stories. To memorize each PAO story the user would be shown four images: a person, an action, an object and a scene. The user is instructed to imagine the PAO story taking place inside the scene. After the user has memorized a PAO story the computer stores the images of the person and the scene, but discards the images of the action and object. The images of the person and the scene are used as a public cue to help the user remember the secret action and object. A password is formed by concatenating the secret action(s) and object(s) from several different PAO stories. The images of the corresponding people/scenes are used as a public cue to help the user remember his secret stories. We stress that the actions and the objects in each of these stories are randomly chosen by the computer after the images of person/scene have been fixed. If the user selected the action and the the object then he might pick actions or objects that are correlated with the person or the scene (e.g., the user might pick the object 'apple' with Steve Jobs). By having the computer select the story we ensure that the secret actions and objects are not correlated with the public cue for the password.

*Sharing Stories.* Stories are shared across different accounts to minimize the total number of stories that the user needs to remember and, more importantly, to maximize natural rehearsals for each PAO story from the user's normal login habits. A central idea in this construction to balance usability and security is $(n, \ell, \gamma)$-sharing set families.

*Definition 1:* We say that a set family $\mathcal{S} = \{S_1, ..., S_m\}$ is $(n, \ell, \gamma)$-sharing if (1) $|\bigcup_{i=1}^{m} S_i| = n$, (2) $|S_i| = \ell$ for each $S_i \in \mathcal{S}$, and (3) $|S_i \cap S_j| \leq \gamma$ for each pair $S_i \neq S_j \in \mathcal{S}$.

Here, $n$ denotes the number of secrets that the user has to memorize and $m$ denotes the number of passwords that the user can form (e.g., the password $pw_i$ for account $A_i$ is formed by appending all $\ell$ secrets from set $S_i$ together). If the adversary learns all of the secrets in the set $S_j$ (e.g., in a plaintext password breach) then the password for account $S_i$ is still at least as strong as a password containing $|S_i| - \gamma$ secrets.

Blocki et al. [25] showed how to generate $m = 110$ passwords from 43 PAO stories using a $(43, 4, 1)$-sharing set family (each of the 43 secrets is a random action-object pair). In their scheme each password consists of a subset of four random action-object pairs. This scheme provides strong security guarantees. Even after two plaintext password breaches

each of the remaining passwords is strong enough to resist an offline attack assuming that the password was encrypted by a password hash function $\mathbb{H}$ with $\mathbb{C}(\mathbb{H}) \geq \$10^{-6}$ and the adversary is not willing to spend more than $30 cracking the password [24]. Even after three plaintext password breaches the remaining passwords are still strong enough to resist online attacks[3].

*Usability.* To ensure that the user remembers all of his secret PAO stories, he is reminded to rehearse each PAO story that he has not used recently enough. More formally, given a constant $\sigma_s > 1$, which may depend on the strength of the mnemonic techniques used to memorize the secret $s$, and a base unit of time $b$ the user is reminded to rehearse at time $b\sigma_s^{i+1}$ whenever the user has not naturally rehearsed (e.g., by authenticating at an account whose password involves the secret $s$) the secret $s$ during days $[b\sigma_s^i, b\sigma_s^{i+1})$. The usability of a password management scheme is evaluated by predicting how many extra rehearsals ($XR_\infty$) the user would need to perform over his lifetime to remember all of his secrets.

### C. Variants Considered in Our Study

We observe that each PAO story that the user memorizes in the Shared Cues scheme could be viewed as containing two secrets (the action and the object). In this case a user who has memorized four PAO stories would be able to create $m = 14$ secure passwords by adopting the Shared Cues scheme with the following $(8, 4, 2)$-sharing set family $\mathcal{S} = \{\{1, 2, 3, 4\}, \{1, 2, 5, 6\}, \{1, 2, 7, 8\}, \{1, 3, 5, 7\}, \{1, 3, 6, 8\}, \{1, 4, 5, 8\}, \{1, 4, 6, 7\}, \{2, 3, 5, 6\}, \{2, 3, 6, 7\}, \{2, 4, 5, 7\}, \{2, 4, 6, 8\}, \{3, 4, 5, 6\}, \{3, 4, 7, 8\}, \{5, 6, 7, 8\}\}$.

*Security.* Each password is strong enough to resist an offline attack assuming that the password was encrypted by a password hash function $\mathbb{H}$ with $\mathbb{C}(\mathbb{H}) \geq \$10^{-6}$ and the adversary is not willing to spend more than $30 cracking the password [24]. Even if the adversary recovered one of the passwords in a plaintext password breach all of the user's other passwords will most likely be safe against online attacks because each password will contain at least two unknown secrets (action(s) and/or object(s))[4].

*Usability.* The evaluation of usability of this construction can be decomposed into two questions. First, can users robustly recall 4 PAO stories while following a suitable spaced repetition schedule? A central goal of our study is to answer this question. Second, how many extra rehearsals (beyond rehearsals from normal logins) does a user have to perform in order to follow the spaced repetition schedule? We do not attempt to answer this question in our study. However, we provide a sense of this user effort in the discussion section (see Section VI).

[3]We assume that a $k$-strikes policy is used to limit the number of incorrect guesses at each account and that the value of $k$ is reasonably small (e.g., 3, 5)

[4]We remark that there is a $(8, 4, 3)$-sharing set family of size $m = 70 = \binom{8}{4}$. However, we could not create 70 secure passwords using this set family — an adversary who has seen one of the user's plaintext passwords could most likely crack at least one of the user's other passwords in an online attack even with a 3 strikes policy at each account.

## III. Study Design

Our user study was conducted online using Amazon's Mechanical Turk framework, on a website at our institution. It was approved by the Institutional Review Board (IRB) at Carnegie Mellon University under IRB protocol HS14-294: Sufficient Rehearsal Schedules and Mnemonic Techniques. After participants consented to participate in the research study, we randomly assigned each participant to a particular study condition. Members in a particular condition were assigned a particular number of action-object pairs (either 1, 2, or 4), a particular memorization technique (e.g., mnemonic or text), and a particular rehearsal schedule (e.g., 24hr×2, 12hr×1.5) as determined by the condition.

Participants were then asked to complete a memorization phase. We randomly selected actions (e.g., swallowing) and objects (e.g., bike) for each participant to memorize. Participants in mnemonic conditions were also assigned pictures or "scenes," one for each action-object pair, and were given specific instructions about how to memorize their words. We paid participants $0.50 for completing the memorization phase. Once participants completed the memorization phase we asked them to return periodically to rehearse their words. To encourage participants to return we paid participants $0.75 for each rehearsal, whether or not they were able to remember the words. If a participant forgot an action-object pair, then we reminded the participant of the actions and objects that were assigned and asked that participant to complete the memorization phase again.

We restricted our participant pool to those Mechanical Turk workers who had an approval rate of 95% or better, had completed at least 100 previous tasks, and were identified by Amazon as living in the United States. 797 participants visited our study website, and 578 completed the memorization phase and initial rehearsal phase.

### A. Recruitment

On the Mechanical Turk website, participants were recruited with the following text:

> Participate in a Carnegie Mellon University research study on memory. You will be asked to memorize and rehearse random words for a 50 cent payment. After you complete the memorization phase, we will periodically ask you to return to check if you still remember the words. If you forget the words then we will remind you of the words and ask you to complete the memorization phase again. You will be paid 75 cents upon the completion of each rehearsal.

> Because this is a memory study we ask that you do not write down the words that we ask you to memorize. You will be paid for each completed rehearsal phase — even if you forgot the words.

### B. Memorization Phase

*1) Mnemonic group:* We first describe the memorization phase for participants assigned to a mnemonic condition. Participants in the mnemonic group were given the following instructions:

> This study is being conducted as part of a Carnegie Mellon University research project. It is important that you answer questions honestly and completely. Please take a minute to read the following instructions.

> The goal of this study is to quantify the effects of rehearsal and the use of mnemonic techniques on long term memory retention. In this study you will be asked to memorize and rehearse eight random words (four actions and four objects). During the first phase we will ask you to memorize the eight random words — you will be paid $0.50 upon completion of the memorization phase. After you complete the memorization phase we will periodically ask you to return via email to check if you still remember the words. If you forget the words, we will remind you of the words and ask you to complete the memorization phase again. You will be paid $0.75 upon the completion of each rehearsal.

> **Important:** Because this is a memory study we ask that you do not write down the words we ask you to memorize. You will be paid for each completed rehearsal phase — even if you forgot the words.

> You have been assigned to the mnemonic group, which means that we give you specific instructions about how to memorize the words. One of the purposes of this study is to determine how effective certain mnemonic techniques are during the memorization task. We ask that you follow the directions exactly — even if you would prefer to memorize the words in a different way.

After participants finished reading the instructions the memorization phase proceeded as follows:

**Step 1** Initially, participants were shown a photo of an assigned scene (e.g., Figure 1a). Participants were next asked to select a famous person or character from a predefined list (e.g., Darth Vader) and were shown a photograph of the famous person that they selected — see Figure 1b. Once selected, participants could not change their person choices. After selecting a person, participants saw a randomly selected action-object pair. See Appendix A for the lists of people, actions and objects used in the study.

**Step 2** As shown in Figure 2, we asked participants to imagine a story in which the person they selected is performing the action in the given scene (e.g., imagine Darth Vader bribing the roach on the lily pad). We asked participants to type in this story, with all words in the correct order (Person-Action-Object).

**Step 3** Participants were then required to select photographs of the action and object, and type in the action and object

(a) Scene: Lily Pads on the Amazon River



(b) Person: Darth Vader

Fig. 1: Memorization Step 1. Scene and Person.



Fig. 2: Memorization Steps 2–3. Darth Vader bribing a roach on the lily pad.



Fig. 3: Memorization Steps 1–3 for Text group.

words two more times in separate fields.

We asked most participants to repeat Steps 1 through 3 four times using a new scene (e.g., a baseball field or a hotel room underneath the sea), a new famous person/character and a new action-object pair during each repetition. Thus, most participants memorized a total of eight words (four actions and four objects). After the memorization phase, we asked participants to complete a rehearsal phase (See Figure 4) before leaving the website.

*2) Text group:* We next describe the memorization phase for participants assigned to a text condition. Participants in the text group were given the same instructions as those in the mnemonic group, with the exception of the last paragraph, which begins "You have been assigned to the mnemonic group..." This paragraph is omitted for those in a text condition. After participants finished reading the instructions, the memorization phase proceeded as follows:

**Step 1** As shown in Figure 3, we randomly selected an action-object pair, and displayed these words to the participant.

**Step 2** We asked each participant to spend one minute memorizing his words. We suggest that participants imagine a person performing the action with the object. We asked each participant to type in a story which includes the action and the object in the correct order.

**Step 3** Participants were then required to type in the action and object words two more times in separate fields.

As with the mnemonic group, we asked most participants to repeat Steps 1 through 3 four times, and asked participants to complete a rehearsal phase.

*C. Rehearsal Phase*

Each participant was assigned a particular rehearsal schedule. The particular times that we ask the participant to return were given by the rehearsal schedule that participant was assigned to use (see Table II). We e-mailed participants to remind them to return for each rehearsal using the following text:

Dear Carnegie Mellon study participant: Please return to (url) to participate in the next part of the
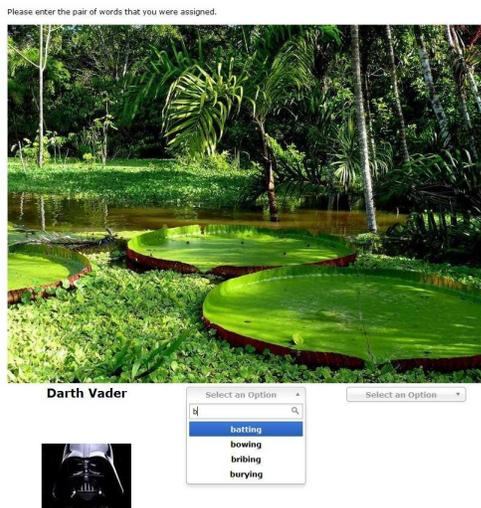
Fig. 4: Rehearsal Phase. Darth Vader and the photo of the lily pads on the Amazon River provide a cue to aid memory recall.

memory study. If you do not return promptly upon receiving this email, you might not be considered for future phases of the study. You will receive a $0.75 bonus payment for completing this task and it should take less than five minutes.

Remember that you should not write down the words that were assigned to you. You will be paid for each completed rehearsal phase — even if you forgot the words.

There is no need to return to Mechanical Turk and find the HIT to receive the bonus, this bonus and any future bonuses will be applied to this MTurk account automatically as you complete each phase. Please do not attempt to take the HIT again on MTurk as this will result in a rejection.

If, for any reason, you do not want to complete the study, please reply to this email and let us know why, so we can improve our protocol for future studies.

Thank you! The Carnegie Mellon University Study Team

We describe the rehearsal phase below:

*1) Mnemonic group:* Each participant was shown the scene and the picture of the person that he chose while memorizing his first story during the memorization phase. We then asked each participant to recall the assigned action and object for that story. As shown in Figure 4, actions and objects were browsable and searchable to aid recall. If the participant was correct then we moved on to the next story. If the participant was incorrect then we asked the participant to try again. Each participant was allowed three guesses per action-object pair. After three incorrect guesses we asked the participant to repeat the memorization phase with the same actions and objects, and try the rehearsal phase again. Once the participant correctly entered all assigned action-object pairs, the rehearsal phase

ended and participants were paid automatically.

*2) Text group:* Each participant from the text group was simply asked to recall the actions and objects assigned during the memorization phase. As with the mnemonic group, actions and objects were browsable and searchable to aid recall. Scoring and payment were handled the same as in the mnenomic conditions.

### D. Follow Up Survey

Some participants did not return to rehearse their stories during the rehearsal phase. We cannot tell whether or not these participants would have remembered their passwords if they had returned. Instead we can only report the fraction of participants who remembered their passwords among those who returned for each rehearsal during the study. There are several reasons why a participant may not have returned (e.g., too busy, did not get the follow up message in time, convinced he or she would not remember the password). If participants do not return because they are convinced that they would not remember the password then this could be a source of bias (i.e., we would be selecting participants who are confident that they remember the story). Our hypothesis is that the primary reason that participants do not return is because they are too busy, because they did not get our follow up message in time, or because they are not interested in interacting with us outside of the initial Mechanical Turk task, and not because they were convinced that they would not remember the story. In order to test our hypothesis we sent a follow up survey to all participants who did not return to complete a rehearsal phase. Participants were paid 25 cents for completing this survey. The survey is described below:

You are receiving this message because you recently participated in a CUPS Memory Study at CMU. A while ago you received an e-mail to participate in a follow up test. We would like to ask you you to complete a quick survey to help us determine why participants were not able to return to complete this follow up study. The survey should take less than a minute to complete, and you will be paid 25 cents for completing the survey. The survey consists of one question. Which of the following reasons best describes why you were unable to return to take the follow up test?

A I no longer wished to participate in the study.
B I was too busy when I got the e-mail for the follow up test.
C I did not see the e-mail for the follow up test until it was too late.
D I was convinced that I would not be able to remember the words/stories that I memorized when I received the e-mail for the follow up test.
E I generally do not participate in follow up studies on mechanical turk.

It is possible that some participants will choose not to participate in the follow up survey. However, in our case their

decision not to participate is valuable information which supports our hypothesis, i.e., they are not interested in interacting with us outside of the initial Mechanical Turk task.

### E. Rehearsal Schedules

In our study each participant was randomly assigned to follow one of the following rehearsal schedules 24hr×2, 24hr×2+2start, 30min×2 and 12hr×1.5. The specific rehearsal times for each schedule can be found in Table II. We can interpret a schedule 24hr×2 as follows: the length of the first rehearsal interval (e.g., the time between initial visit and the first rehearsal) is 24 hours and the length of the $i + 1$'th rehearsal interval is twice the length of the $i$'th rehearsal interval. If the participant was assigned to the 24hr×2 rehearsal schedule then we would send that participant a reminder to rehearse 1 day after the memorization phase. If that participant successfully completes the first rehearsal phase then we will send that participant another reminder to rehearse 2 days after the first rehearsal, and the next reminder would come four days later, etc. The final rehearsal would take place on day $1 + 2 + ... + 32 + 64 = 127$. In the 12hr×1.5 schedule the length of the first rehearsal interval is 12 hours and after that intervals grow by a factor of 1.5. The 24hr×2+2start and 30min×2 conditions are similar to the 24hr×2 schedule except that participants are asked to do a few additional rehearsals on day 1 — after this the rehearsal intervals are identical to the 24hr×2 schedule.

We use the following syntactic pattern to denote a study condition: (Memorization Technique)_(Rehearsal Schedule)_(Number of action-object pairs memorized). For memorization technique we use $m$ to denote the mnemonic groups and $t$ to denote the text group. Thus, a participant in the group m_24hr×2_4 refers to a user who was asked to memorize four actions and four objects using the mnemonic techniques we suggested and to rehearse his person-action-object stories following the 24hr×2 rehearsal schedule from Table II. Because most participants were asked to memorize four random action-object pairs we drop the "_4" from the end of those conditions.

### F. Online studies

The passwords in our study did not protect high-value accounts, limiting ecological validity. In contrast to real-world, high-value passwords, study participants would not suffer consequences beyond a modest time cost if they forgot their password, nor were they incentivized to keep their passwords only in memory beyond our repeated requests that they do so.

We recruited participants using Mechanical Turk (MTurk). Using MTurk allows us to study a larger volume of participants in a controlled setting than would otherwise be possible. MTurk workers tend to be younger, more educated, and more technical than the general population, but they represent a significantly more diverse population than is typically used in lab studies, which often rely on college-student participants [26], [27]. Many researchers have found that well-designed MTurk studies provide high-quality user data [28]–[33]. Adar has criticized MTurk studies in general, although

our use of crowdsourcing to understand human behavior fits his description of an appropriate use [34].

## IV. RESULTS

In this section we present the results from our study. In Section IV-A we overview the raw data from our study (e.g., how many participants returned for each rehearsal round?) and some simple metrics (e.g., how many of these participants remembered their action-object pairs?) In Section IV-B we discuss the results of a survey we sent to participants that did not return for a rehearsal phase. In Section IV-C we briefly overview Cox regression — a tool for performing survival analysis that we used to compare several of our study conditions. In Section IV-D we use the data from our study to evaluate and compare different study conditions.

### A. Study Data

We first overview the metrics we use to evaluate the performance of participants in different study conditions.

**Notation:** Given a participant $P$ we use the indicator function **Returned** $(P, i) = 1$ (resp. **Remembered** $(P, i) = 1$) if and only if $P$ returned for rehearsal $i$ (resp. if and only if $P$ remembered his words during rehearsal $i$ with $< 3$ incorrect guesses per action-object pair. ). We use the function **Survived** $(P, i) = \prod_{j=1}^{i}$ **Remembered** $(P, i)$ to indicate whether $P$ remembered his words with $< 3$ incorrect guesses per action-object pair during rehearsal $i$ and during every earlier rehearsal $j < i$. We use the function **SuccessfulReturned** $(P, i) =$ **Survived** $(P, i - 1) \wedge$ **Returned** $(P, i)$ to indicate whether $P$ survived rehearsals 1 to $i - 1$ with no failures and returned for rehearsal $i$. Given an indicator function **F** and a study condition $C$ (e.g., a set of participants) we use

$$\mathbf{NumF}(C, i) = \sum_{P \in C} \mathbf{F}(P, i)$$

to denote the number of participants selected by the indicator function **F**. For example, **NumSurvived** $(C, i)$ denotes the number of participants in condition $C$ who survive through rehearsal $i$. We will usually omit the $C$ and write **NumSurvived** $(i)$ when discussing results within a particular condition. Finally, we use **Time** $(i)$ to denote the time of rehearsal $i$, as measured from the initial memorization phase.

Table III shows how many participants who had never failed before returned in each rehearsal round as well as their conditional probability of success with 95% confidence intervals (e.g., **NumSuccessfulReturned** $(i)$ and **NumSurvived** $(i)$ /**NumSuccessfulReturned** $(i)$). We note that in some study conditions there is still one rehearsal round that has not been completed yet — these conditions are still ongoing. Figures 5a and 5b plot the conditional probability of success for participants who have not failed before (e.g., **NumSurvived** $(i)$ /**NumSuccessfulReturned** $(i)$). Figure 6 plots the probability of success for all

| Schedule | Multiplier | Base | Rehearsal Intervals | Rehearsal Days |
|---|---|---|---|---|
| 24hr×2 | ×2 | 1 Day | 1, 2, 4, 8, 16, 32, 64 | 1, 3, 7, 15, 31, 63, 127 |
| 12hr×1.5 | ×1.5 | 0.5 days | 0.5, 1.25, 2.4, 4, 6.5, 10,16,24,37,56 | 0.5, 1.75, 4.15, 8.15, 14.65, 24.65, 40.65, 64.65, 101.65, 157.65 |
| 24hr×2+2start | ×2 | 1 Day | 0.1 days, 0.5, 1, 2, 4, 8, 16, 32, 64 | 0.1, 0.6, 1.6, 3.6, 7.6, 15.6, 31.6, 63.6, 127.6 |
| 30min×2 | ×2 | 30 min | 0.5 hr, 1 hr, 2 hr, 4 hr, 8 hr, 1 day, 2 , 4 , 8, 16 , 32 , 64 | 0.5hr, 1.5hr, 3.5hr, 7.5hr, 15.5hr, 1.65 days, 3.65, 7.65, 15.65, 31.65, 63.65, 127.65 |

TABLE II: Rehearsal Schedules

| Rehearsal $i$\ Condition | Initial | NumSuccessfulReturned$(i)$, NumSurvived$(i)$/NumSuccessfulReturned$(i)$, 95% confidence interval | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $i = 0$ | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| m_24hr×2+2start | 80 | 51 94.1% 0.838,0.988 | 41 100% 0.914,1 | 38 100% 0.907,1 | 37 97.3% 0.858,0.999 | 36 100% 0.903,1 | 36 97.2% 0.855,0.999 | 34 97.1% 0.847,0.999 | 30 90% 0.735,0.978 | 24 68% 0.465,0.85 |
| t_24hr×2+2start | 100 | 71 88.7% 0.790,0.950 | 54 96.3% 0.873,0.995 | 51 100% 0.930,1 | 51 100% 0.930,1 | 51 94.1% 0.836,0.988 | 48 95.8% 0.857,0.995 | 44 88.6% 0.754,0.962 | 39 79.5% 0.635,0.907 | 39 86.2% 0.683,0.961 |
| m_24hr×2 | 75 | 65 76.9% 0.648,0.845 | 45 93.3% 0.817,0.986 | 40 100% 0.912,1 | 39 97.4% 0.865,0.999 | 37 97.3% 0.858,0.999 | 32 96.9% 0.838,0.999 | | | |
| m_24hr×2+2start_2 | 81 | 50 100% 0.929,1 | 42 100% 0.916,1 | 42 100% 0.916,1 | 41 100% 0.914,1 | 38 100% 0.907,1 | 37 100% 0.905,1 | 36 100% 0.903,1 | 33 100% 0.894,1 | |
| m_24hr×2+2start_1 | 86 | 64 100% 0.943,1 | 52 100% 0.932,1 | 49 100% 0.927,1 | 49 100% 0.927,1 | 47 100% 0.925,1 | 46 100% 0.923,1 | 45 100% 0.922,1 | 44 100% 0.920,1 | |
| m_12hr×1.5 | 83 | 72 86.1% 0.759,0.931 | 53 98.1% 0.899,1.000 | 51 100% 0.930,1 | 51 100% 0.930,1 | 49 100% 0.927,1 | 46 97.8% 0.885,0.999 | 43 100% 0.918,1 | 42 97.6% 0.874,0.999 | 42 94.9% 0.827,0.994 |
| m_30min×2 | 73 | 40 95% 0.831,0.994 | 27 100% 0.872,1 | 26 100% 0.868,1 | 24 100% 0.858,1 | 22 100% 0.846,1 | 22 100% 0.846,1 | 22 100% 0.846,1 | 22 100% 0.846,1 | 22 100% 0.846,1 |
| Rehearsal $i$\ Condition | | $i = 10$ | $i = 11$ | | | | | | | |
| m_30min×2 | | 21 95.2% 0.762,0.999 | 20 90% 0.683,0.988 | | | | | | | |

TABLE III: **NumSurvived**$(i)$/**NumSuccessfulReturned**$(i)$ with 95% binomial confidence intervals. $m =$ "mnemonic," $t =$"text"

participants who returned for rehearsal $i$ (e.g., **NumRemembered**$(i)$ /**NumReturned**$(i)$) with corresponding values in Table IV.

One of the primary challenges in analyzing the results from our study is that some participants were dropped from the study because they were unable to return for one of their rehearsals in a timely manner. We do not know how these participants would have performed under ideal circumstances (e.g., if all of our participants were always able to return to rehearse in a timely manner). We compare three different metrics to estimate the survival rate of participants in our study under ideal circumstances.

Our first estimate is **EstimatedSurvival**$(i) =$

$$\prod_{j=1}^{i} \frac{\textbf{NumSurvived}\,(j)}{\textbf{NumSuccessfulReturned}(j)} \,,$$

where $\frac{\textbf{NumSurvived}(j)}{\textbf{NumSuccessfulReturned}(j)}$ denotes our empirical estimate of the conditional probability that a participant will survive round $j$ given that the participant survived all previous rounds and returned for rehearsal $j$. We plot this value in Figures 7a and 7b.

Our second estimate, shown in Figure 8, is much simpler:

$$\textbf{ObservedSurvival}\,(i) = \frac{\textbf{NumSurvived}\,(i)}{\textbf{NumReturned}\,(i)} \,.$$

One potential concern with this estimate is that participants who failed in earlier rehearsal rounds might be less likely to return for future rehearsals. If so, this would bias the estimate upward, as later rounds would have a biased sample of participants who survived previous rehearsals. However, we did not observe any obvious correlation between prior failure and the return rate. Sometimes the return rate was higher for participants who had failed earlier than for participants who had never failed and sometimes the return rate was lower. Furthermore, in our survey of participants who did not return in time for a rehearsal round no one self-reported that they did not return because they were not confident that they would be able to remember (see Section IV-B). Both our first and second estiamtes were consistently close.

| Rehearsal $i\backslash$ Condition | Initial | Returned $(i)$, NumRemembered$(i)$/NumReturned$(i)$, 95% confidence interval | | | |
|---|---|---|---|---|---|
| $i=0$ | | 6 | 7 | 8 | 9 |
| m_24hr×2+2start | 80 | 40 87.5% 0.732,0.958 | 38 86.8% 0.719,0.956 | 34 79.4% 0.621,0.913 | 31 54.8% 0.360,0.727 |
| t_24hr×2+2start | 100 | 58 79.3% 0.666,0.888 | 56 69.6% 0.559,0.812 | 55 56.4% 0.423,0.697 | 50 50% 0.355,0.645 |
| m_24hr×2 | 75 | 42 73.8% 0.580,0.861 | | | |
| m_24hr×2+2start_2 | 81 | 37 100% 0.905,1 | 36 100% 0.903,1 | 33 100% 0.894,1 | |
| m_24hr×2+2start_1 | 86 | 46 100% 0.923,1 | 45 100% 0.922,1 | 44 100% 0.920,1 | |
| m_12hr×1.5 | 83 | 55 81.8% 0.691,0.909 | 53 81.1% 0.680,0.906 | 51 80.4% 0.669,0.902 | 48 77.1% 0.627,0.880 |
| m_30min×2 | 73 | 22 100% 0.846,1 | 22 100% 0.846,1 | 22 100% 0.846,1 | 22 100% 0.846,1 |

| Rehearsal $i\backslash$ Condition | | $i=10$ | $i=11$ | | |
|---|---|---|---|---|---|
| m_30min×2 | | 21 95.2% 0.762,0.999 | 21 85.7% 0.637,0.970 | | |

TABLE IV: **NumRemembered**$(i)$/**NumReturned**$(i)$ with 95% binomial confidence intervals. $m =$ "mnemonic," $t =$"text"

We also include a pessimistic estimate of the survival rate, i.e., an estimate that is biased downward. Figure 9 plots the value of **PessimisticSurvivalEstimate**$(i) =$

$$\frac{\textbf{NumSurvived}(i)}{\textbf{NumSuccessfulReturned}(i) + \sum_{j=1}^{i}\textbf{NumFailed}(j)}$$

where the indicator function **Failed**$(P, j) = 1$ if and only if **Survived**$(P, j - 1)$ and $P$ failed to remember at least one of his action-object pairs with $< 3$ guesses during rehearsal $j$. This estimate explicitly assumes that participants who failed previously chose not to return, and is most likely overly pessimistic. A participant who succeeded through rehearsal $i$ but could not return for rehearsal $i+1$ would not be included in the estimate. However, if our participant had failed during round $i$ before not returning for rehearsal $i+1$ then s/he would still be included as failing at $i + 1$.

### B. Survey Results

We surveyed 61 participants who did not return to complete their first rehearsal to ask them why they were not able to return. The results from our survey are presented in Figure 10. The results from our survey strongly support our hypothesis that the primary reason that participants do not return is because they were too busy, because they did not get our follow up message in time, or because they were not interested in interacting with us outside of the initial Mechanical Turk task, and not because they were convinced that they would not remember the story — no participant indicated that they did



(a) Faceted. Mean Time Since Memorization.



(b) Together. Mean Time Since Memorization.

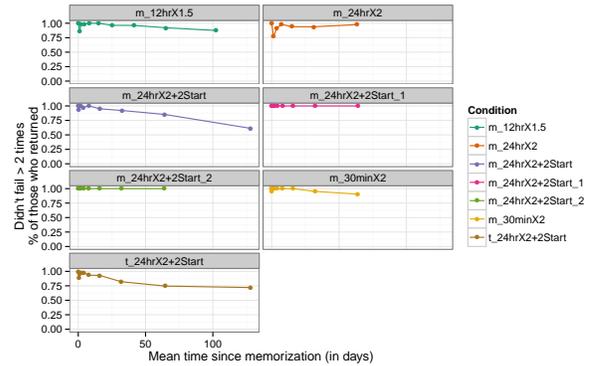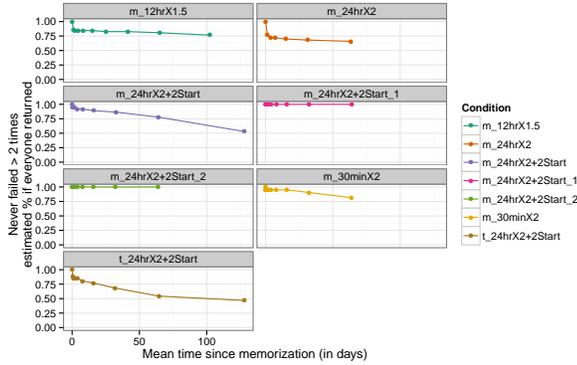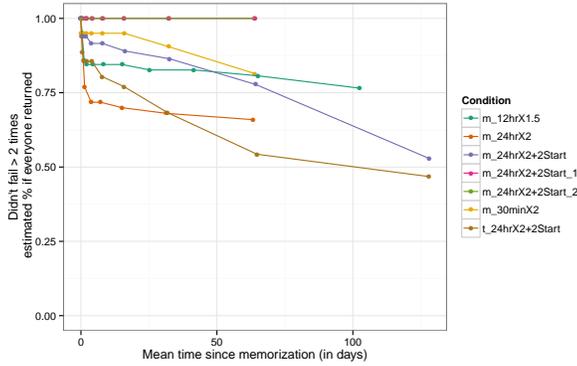Fig. 5: **NumSurvived**$(i)$/**NumSuccessfulReturned**$(i)$ vs. **Time**$(i)$



Fig. 6: **NumRemembered**$(i)$/**NumReturned**$(i)$ vs **Time**$(i)$

not return because they thought that they would not be able to remember the action-object pairs that they memorized.

*Fun:* We had several participants e-mail us to tell us how much fun they were having memorizing person-action-object stories. The results from our survey are also consistent with the hypothesis that memorizing person-action-object stories is fun (e.g., no participants said that they no longer wished to participate in the study).

(a) Faceted.



(b) Together.

Fig. 7: $\mathbf{EstimatedSurvival}\,(i)$ vs $\mathbf{Time}\,(i)$
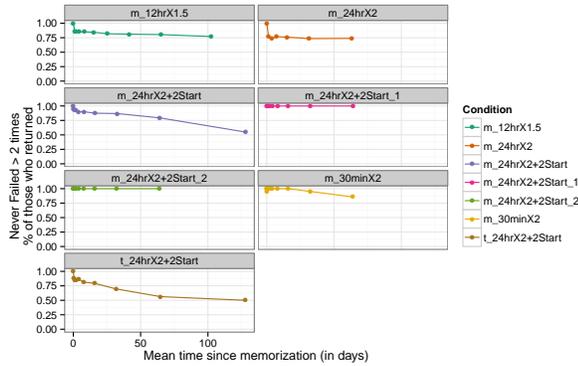


Fig. 8: $\mathbf{ObservedSurvival}\,(i)$ vs $\mathbf{Time}\,(i)$

*C. Statistical methods*

We used Cox regression to perform survival analysis, and compare different study conditions. Survival analysis relates the time that passes before a failure event (e.g., a user forgets one of his action-object pairs during a rehearsal) to covariates (e.g., mnemonic/text, rehearsal schedule, number of words memorized) that may be associated with this quantity of time. Cox regression is an appropriate tool for our study because it can deal with participants who dropped out of the study before they failed (e.g., participants who did not return for a rehearsal in a timely manner).



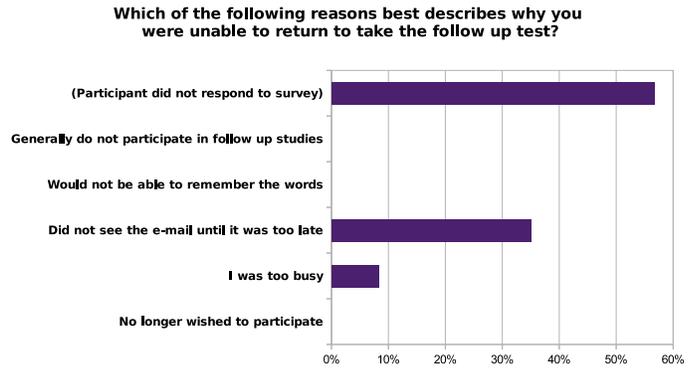Fig. 9: $\mathbf{PessimisticSurvivalEstimate}\,(i)$ vs $\mathbf{Time}\,(i)$.



Fig. 10: Survey Results

To run Cox regression we need to have the following data for each participant: whether or not they failed, the time of failure for participants who did fail and the time a participant dropped out of the study for participants who were dropped before their first failure. The output of Cox regression is a set of regression coefficients $\beta_1, \ldots, \beta_k$, which tell us how different study conditions affect the survival rate[5]. For example, suppose that our baseline condition was the t_24hr$\times$2+2start condition and that $x_1 = 1$ if and only if the user was in the m_24hr$\times$2+2start condition. If $\exp(\beta_1) < 1$ (resp. $\exp((\beta_1)) > 1$) then participants in the mnemonic condition are less likely (resp. more likely) to fail at any given time than participants in the baseline text condition.

We contrast the results of Cox regression, which operates over the full duration of our study, with simple $t$-tests performed on estimated survival rates between conditions at a

[5] We use the proportional hazard model to compare the risk of failure for participants in different study conditions. In particular, given a baseline study condition we let $\lambda_0(t)$ denote the underlying hazard function — a function describing the risk that a participant fails to remember his action-object pairs at time $t$. Given explanatory variables $x_1, \ldots, x_p$ we use the function

$$\lambda\,(t \mid x_1, \ldots, x_p) = \lambda_0\,(t) \exp\,(\beta_1 x_1 + \ldots + \beta_p x_p)\,,$$

to compare the risk for participants in different study conditions and we use Cox regression to compute the regression coefficients $\beta_1, \ldots, \beta_p$ for each explanatory variable.

| Condition $x_i$ | $\beta_i$ | $\exp(\beta_i)$ | $\exp(-\beta_i)$ | 95% Confidence Interval for $\exp(\beta_i)$ | |
|---|---|---|---|---|---|
| | | | | Lower | Upper |
| m_30min×2 | -0.7334 | 0.4803 | 2.082 | 0.1793 | 1.2867 |
| m_24hr×2+2start | -0.9645 | 0.3812 | 2.624 | 0.1658 | 0.8762 * |
| m_12hr×1.5 | -0.7454 | 0.4746 | 2.107 | 0.2366 | 0.9519 * |

$n = 228$, number of failure events $k = 56$.
\* indicates $\exp(\beta_i)$ is significantly different from 1 at the $\alpha = 0.05$ level.

TABLE V: Cox regression with baseline: m_24hr×2

| Condition $x_i$ | $\beta_i$ | $\exp(\beta_i)$ | $\exp(-\beta_i)$ | 95% Confidence Interval for $\exp(\beta_i)$ | |
|---|---|---|---|---|---|
| | | | | Lower | Upper |
| t_24hr×2+2start | 0.4183 | 1.519 | 0.6582 | 0.843 | 2.738 |

$n = 122$, number of failure events $k = 49$.

TABLE VI: Baseline: m_24hr×2+2start

specific point in time, typically 64 days after memorization. We do this to examine estimated performance after passwords have been used for an extended period of time.

### D. Findings

*1) Reliable Recall of Multiple Passwords:* We find that 77.1% of participants in the m_12hr×1.5 condition who were always able to return were able to remember all four of their secret action-object pairs during all nine of their rehearsal rounds of a time period of 101.65 days. This value ($\mathbf{NumSurvived}(9)/\mathbf{NumReturned}(9) = 0.771$) closely matches our estimated survival rate $\mathbf{EstimatedSurival}(9) = 0.765$. As explained before, eight secrets (four actions and four objects) could be used to form 14 different passwords using the Shared Cues password management scheme [12]. As we can see from Table III, most of the participants who did not survive forgot their action-object pairs during the first rehearsal.

*2) Reliable Recall of a Single Password:* 100% of participants who were asked to memorize one or two PAO stories (e.g., one password) were able to remember all of their secret action-object pairs during all eight rehearsals over a time period of 63.6 days. The 95% confidence interval for the fraction of participants in the m_24hr×2+2start_1 (resp. m_24hr×2+2start_2) condition who always remember their action-object pairs is $[0.920, 1]$ (resp. $[0.894, 1]$).

*3) Effect of Rehearsal:* We used Cox regression to compare the survival rate for participants in the m_24hr×2, m_12hr×1.5, m_24hr×2+2start, and m_30min×2 study conditions. We used the m_24hr×2 condition as our baseline. Our results are shown in Table V. For all three conditions m_12hr×1.5, m_24hr×2+2start, and m_30min×2 we have $\exp(\beta_i) < 1$ meaning that our model predicts that participants in the baseline condition (m_24hr×2) were less likely to survive at any given point in time. The evidence for the hypothesis $\exp(\beta_i) < 1$ is statistically significant (at the $\alpha = 0.05$ level) for the m_12hr×1.5 and m_24hr×2+2start conditions because the confidence interval for $\exp(\beta_i)$ does not contain the value 1. We cannot claim statistical significance for the 30min×2 condition due to a low number of participants in this group[6].

We also used used one-tailed $t$-tests to test the hypothesis $\mathbf{ObservedSurvival}(C_1, i_1) >$

[6]Because the 30min×2 condition had many rehearsals on day 1 our time window for participants to return for each of these rehearsals was more narrow than in other conditions. As a result many participants in this condition were dropped after day 1 because they were not able to return in a timely manner.

$\mathbf{ObservedSurvival}(C_2, i_2)$ for $C_1 = $ m_24hr×2 and $C_2 \in \{$ m_12hr×1.5, m_30min×2, m_24hr×2+2start $\}$. We adjusted the values of $i_1$ and $i_2$ to compare the survival rates over similar time periods for participants who followed different rehearsal schedules. In particular, every rehearsal schedule had one rehearsal near day 64 (e.g., the sixth rehearsal in the 24hr×2 schedule was on day 63 and the eighth rehearsal for the 12hr×1.5 schedule was on day 64.65). Table VIII shows these results. While the survival rates were lowest in the m_24hr×2 condition the $t$-test results were not statistically significant at the $\alpha = 0.05$ level.

*4) Mnemonic vs Text Conditions:* We found that participants who used the PAO mnemonic technique significantly outperform users who didn't in recalling their action-object pairs in the short term. In particular, we used a one-tailed $t$-test and tested the hypothesis $\mathbf{ObservedSurvival}(C_1, i_1) > \mathbf{ObservedSurvival}(C_2, i_2)$ for $C_1 = $ m_24hr×2+2start and $C_2 = $ t_24hr×2+2start. We compared our conditions at two points in time: 63.6 days after memorization ($i_1 = i_2 = 8$) and after the final rehearsal 127.6 days after memorization ($i_1 = i_2 = 9$). Table VIII shows these results.

The evidence that participants in m_24hr×2+2start perform better than participants in t_24hr×2+2start for the first 63.6 days after memorization is statistically significant ($p = 0.010$). However, after the final rehearsal ($i_1 = i_2 = 9$) the evidence for our hypothesis is no longer statistically significant. Surprisingly, during the last rehearsal round on day 127.6, participants in the t_24hr×2+2start condition were more likely to remember their action-object pairs than participants in the m_24hr×2+2start condition. We also used Cox regression to compare the survival rates for the t_24hr×2+2start and m_24hr×2+2start conditions using the m_24hr×2+2start condition as a baseline. Our results are shown in Table VI. We have $\exp(\beta_i) > 1$ for the t_24hr×2+2start condition, which indicates that participants did benefit from adopting mnemonic techniques to memorize their action-object pairs. However, the hypothesis $\exp(\beta_1) > 1$ (e.g., the survival rate is worse in the text condition) does not reach the level of statistical significance at the $\alpha = 0.05$ level. Table VII shows the results of Cox regression if we only include data from the first eight rehearsal rounds (through day 63.6). In this case the hypothesis $\exp(\beta_i) > 1$ is statistically significant.

*5) Effect of Interference:* We found that there is an interference effect. Participants performed better when memorizing one or two PAO stories. We used a one-tailed $t$-test to test the hypothesis $\mathbf{ObservedSurvival}(C_1, i_1) > \mathbf{ObservedSurvival}(C_2, i_2)$ for the conditions $C_1 \in \{$m_24hr×2+2start_1, m_24hr×2+2start_2$\}$ and $C_2 = $ m_24hr×2+2start_4. We tested the hypothesis 63.6 days after

| Condition $x_i$ | $\beta_i$ | $\exp(\beta_i)$ | $\exp(-\beta_i)$ | 95% Confidence Interval for $\exp(\beta_i)$ | |
|---|---|---|---|---|---|
| | | | | Lower | Upper |
| t_24hr×2+2start | 1.04 | 2.829 | 0.3434 | 1.287 | 6.219 * |

$n = 122$, number of failure events $k = 37$.
* indicates $\exp(\beta_i)$ is significantly different from 1 at the $\alpha = 0.05$ level.

TABLE VII: Baseline: m_24hr×2+2start

| Condition $C_1$ | $i_1$ (Day) | Condition $C_2$ | $i_2$ (Day) | $p$-value |
|---|---|---|---|---|
| m_24hr×2+2start | 8 (63.6) | t_24hr×2+2start | 8 (63.6) | 0.010 * |
| m_24hr×2+2start | 9 (127.6) | t_24hr×2+2start | 9 (127.6) | 0.338 |
| m_24hr×2+2start_1 | 8 (63.6) | m_24hr×2+2start Remark 1. | 8 (63.6) | 0.042 * |
| m_24hr×2+2start_2 | 8 (63.6) | m_24hr×2+2start Remark 2. | 8 (63.6) | 0.006 * |
| m_24hr×2+2start_1 | 8 (63.6) | m_24hr×2+2start Remark 3. | 8 (63.6) | 0.0011 * |
| m_24hr×2+2start_2 | 8 (63.6) | m_24hr×2+2start Remark 3. | 8 (63.6) | 0.0011 * |
| m_24hr×2+2start_1 | 8 (63.6) | m_24hr×2+2start Remarks 1 and 3. | 8 (63.6) | 0.03 * |
| m_24hr×2+2start_2 | 8 (63.6) | m_24hr×2+2start Remarks 2 and 3. | 8 (63.6) | 0.046 * |
| m_24hr×2+2start | 8 (63.6) | m_24hr×2 | 6 (63) | 0.287 |
| m_30min×2 | 8 (63.6) | m_24hr×2 | 6 (63) | 0.146 |
| m_12hr×1.5 | 8 (64.65) | m_24hr×2 | 6 (63) | 0.228 |

Remark 1. Count participant as surviving if s/he always remembered the first action-object pair.
Remark 2. Count participant as surviving if s/he always remembered the first two action-object pairs.
Remark 3. If a participant dropped and never failed count them as surviving.
* indicates statistical significance at the $\alpha = 0.05$ level.

TABLE VIII: One-Tailed $t$-tests for Hypotheses:
$$\mathbf{ObservedSurvival}(C_1, i_1) >$$
$$\mathbf{ObservedSurvival}(C_2, i_2)$$

memorization (e.g., $i_1 = i_2 = 8$) because this was the last rehearsal round that participants in the m_24hr×2+2start_1 and m_24hr×2+2start_2 conditions completed. The results are shown in Table VIII. The evidence for both hypotheses was statistically significant. In fact, we can confirm a much stronger hypothesis: the survival rate in the m_24hr×2+2start_1 (resp. m_24hr×2+2start_2) condition is greater than the survival rate in the m_24hr×2+2start_4 even if we only count failures on the first (resp. first two) action-object pair(s) and even if we treat every participant $P$ from the m_24hr×2+2start_4 condition who dropped without failing on the first (resp. first two) action-object pair(s) as succeeding. A Cox regression on this data was unable to converge as both $\beta_i$ tended to negative infinity — since both m_24hr×2+2start_1 and m_24hr×2+2start_2 had no failures, $\exp(\beta_i)$ tends to 0.

## V. RELATED WORK

### A. Spaced Repetition

Pimsleur [9] proposed a rehearsal schedule to help people memorize unfamiliar vocabulary words. He suggested rehearsing after 5 seconds, 25 seconds, 2 minutes, 10 minutes, 5 hours, 1 day, 5 days, 20 days, etc. His proposed schedule is precisely the schedule given by expanding rehearsal assumption with the association strength constant set to $2^\sigma = 5 \approx$ and an initial delay before the first rehearsal set to 5 seconds. Pimsleur based his recommendations on previous empirical studies [35, pp. 726 ff]. The application SuperMemo [8] uses a similar rehearsal schedule to help users remember flashcards. Wozniak and Gorzelanczyk conducted an empirical study to test these rehearsal schedules [7]. In their study undergraduate students were asked to memorize and rehearse vocabulary words for a foreign language by following a rehearsal schedule very similar to the expanding rehearsal schedule. Wozniak and Gorzelanczyk tracked each students performance with each particular vocabulary word and used that information to estimate how difficult each word was. If a word was deemed 'difficult' then the length of the time interval before the next rehearsal would only increase by a small multiplicative constant (e.g., 1.5) and if the word was judged to be 'easy' then this time interval would increase by a larger multiplicative constant (e.g., 4).

We stress two key differences in our study: First, because we are asking the user to memorize secrets that will be used to form passwords our rehearsal schedule needs to be conservative enough that our user will consistently be able to remember his secrets during each rehearsal. In other studies the information participants were asked to memorize (e.g., vocabulary words) was not secret so participant could simply look up the correct answer whenever they forgot the correct answer during a rehearsal. By contrast, in the password setting a recovery mechanism may not always be available — users are advised against writing down passwords and organizations have been held liable for damages when they do not properly encrypt their passwords [36]. Second, in our study we ask participants to memorize secrets by following the Person-Action-Object mnemonic techniques. Because these secrets may be easier or harder to memorize than other information like vocabulary words the ideal rehearsal schedule should be tailored to particular mnemonic techniques adopted by the user. Previous studies have demonstrated that cued recall is easier than pure recall (see for example [10]) and that we have a large capacity for visual memories [11]. However, we are not aware of any prior studies which compare cued recall and pure recall when participants are following a rehearsal schedule similar to the one suggested by the expanding rehearsal assumption.

### B. Spaced Repetition – Applications to Passwords

*a) Password Management Schemes:* While there are many articles, books, papers and even comics about selecting strong individual passwords [37]–[44], there has been little work on *password management schemes*—systematic strategies to help users create and remember multiple passwords—that are both usable and secure. Bonneau et al. [1] evaluated several alternatives to text passwords (e.g., graphical passwords, password management software, single-sign-on, federated authentication) finding that, while each alternative had its advantages, none of the alternatives were strictly better than text passwords. Florencio et al. [45] argued that any

usable password management scheme[7] cannot require users to memorize unique random passwords. They suggested that users adopt a tiered password management scheme with a unique password for high, medium and low security accounts. Blocki et al. [12] recently proposed designing password management schemes that maximized the natural rehearsal rate for each of the secrets that the user had to memorize subject to minimum security constraints. Our study is heavily motivated by their work, which we already described in Section II.

*b) Slowly Expanding Password Strength:* Bonneau and Schechter conducted a user study in which participants were encouraged to slowly memorize a strong 56 bit password using spaced repetition [14]. Each time a participant returned to complete a distractor task he was asked to login by entering his password. During the first login the participant was shown four additional random characters and asked to type them in after his password. To encourage participants to memorize these four characters they would intentionally wait a few seconds before displaying them to the user the next time he was asked to login to complete a distractor task. Once a participant was able to login several times in a row (without waiting for the characters to be displayed) they would encourage that participant to memorize four additional random characters in the same way. They found that 88% of participants were able to recall their entire password without any prompting three days after the study was completed. There are several key difference between their study and ours: First, in our study participants were asked to memorize their entire password at the start of the study. By contrast, Bonneau and Schechter encouraged participants to slowly memorize their passwords. Second, Bonneau and Schechter did not tell participants that their goal was to slowly memorize a strong 56 bit password — users were led to believe that the distractor task was the purpose of the study. By contrast, in our study we explicitly told participants that their goal was to remember their words (without writing them down). Finally, participants in our study were given fewer chances to rehearse their passwords and were asked to remember their passwords over a longer duration of time (4 months vs 2 weeks). Bonneau and Schechter asked participants to login 90 times over a two week period. In our study participants were asked to rehearse *at most* 11 times over a period of up to 127 days. We believe that the results of our study could be used to help improve the password strengthening mechanism of Bonneau and Schechter — see discussion in Section VI.

### C. System Assigned Passwords.

Empirical studies have shown that many user-selected passwords are easily guessable [5]. A user study conducted by Shay et al. [46] compared several different methods of generating system assigned passwords for users to memorize (e.g., three to four random words, 5 or 6 random characters). They found that users had difficulty remembering system assigned passwords 48–120 hours after they had memorized it. In fact,

[7]They use the term "password portfolios."

users had more difficulty three to four random words from a small dictionary than when they were asked to remember 5 to 6 random characters. Participants in their study were not asked to follow any particular mnemonic techniques, and were not asked to follow a rehearsal schedule.

### D. Password Composition Policies

Another line of work on passwords has focused on composition policies (policies which restrict the space passwords that users can choose) [40], [47]–[49]. These policies may negatively effect usability (e.g., users report that their passwords are more difficult to remember [40], [47]) and also have adverse security effects (e.g., users are more likely to write down their passwords [47], [49], some restrictive composition policies can actually result in a weaker password distribution [47], [48]).

## VI. Discussion

*Password Expiration Policies:* Following NIST guidelines [13] many organizations require users to change their passwords after certain period of time (e.g., thirty days). The desired behavior is for user's to select a random new password that is uncorrelated with their previous passwords. We contend that these policies will adversely affect usability and security. Memorizing a new password requires effort and users are typically only willing to invest a limited amount time and energy memorizing new passwords. Our experiments indicate that most of the effort to memorize and rehearse a password is spent in the first week after the new password is chosen. By forcing users to reset their password frequently an organization forces its users to remain within the most difficult rehearsal region. There is strong empirical evidence that users respond to password expiration policies by selecting weak passwords and/or selecting new passwords that are highly correlated with one of their old passwords (e.g., old_password+$i$ for $i = 1, 2, \ldots$) effectively canceling out any security gains [50], and many organizations, when faced with competition, avoid password expiration policies entirely due to usability concerns [51]. We contend that a more productive policy would ask participants to slowly strengthen their passwords over time using spaced repetition (see discussion below).

*Strengthening Passwords Over Time:* Our results suggest that the password strengthening mechanism of Bonneau and Schechter [14] could be improved by adopting the PAO story mnemonic used in our study and by using a rehearsal schedule like 24hr×2+2start to help predict when a user has memorized his new secret. Recall that in their mechanism a user authenticates by typing in his old password and then by typing a random character or word that is displayed next to the password box. To encourage participants to memorize this secret character/word, the user will not be shown the random character/word for several seconds, allowing a user who has memorized the secret to authenticate faster who has not. At some point the mechanism will predict that the user has memorized his new secret. At this point this secret is permanently appended to the user's password so that the user must remember this secret to authenticate —

he can no longer wait for the character/word to be displayed. We remark that,instead of requiring the user memorizing a new random character/word to append to his password, it may be easier for the user to memorize a random action-object pair using the PAO mnemonic techniques from this study. Participants in the mnemonic_24hr×2+2start_1 and mnemonic_24hr×2+2start_2 conditions remembered their secret action-object pairs perfectly. We also remark that the 24hr×2+2start rehearsal schedule could provide a reasonable basis for predicting when a user has memorized his new action-object pair. In particular, the 24hr×2+2start schedule can help us predict how long the user will be able to remember his new action-object pairs without rehearsing again. If it is safe to assume that the user will return to authenticate before this point then we would argue that it is safe to predict that the user has memorized his secret action-object pair. Another interesting observation from our study was that participants in the m_24hr×2_4 condition who remembered their action-object pairs during the first two rehearsal on days 1 and 3 were actually more likely to survive through rehearsal 6 (on day 63) than participants in the m_24hr×2+2start_4 condition who remembered their action-object pairs through rehearsal 4 (on day 3.6) were to survive through the corresponding rehearsal 8 (on day 63.6) — though this result was not significant at the $p = 0.05$ level. We hypothesize that a user's ability to remember a particular set of action-object pairs after a challenging rehearsal interval (e.g., only 77% of participants in the aggressive rehearsal who returned for the first rehearsal on day 1 remembered their action-object stories) is better indicator of that user's future success for those particular action-object pairs than performance on less challenging rehearsal intervals. This hypothesis could also help us to predict when a user has memorized a new action object pairs. However, more studies are necessary to properly test this hypothesis.

*Mitigating Initial Forgetting:* Our finding that most participants who did not survive forgot one of their action-object pairs on the first rehearsal (12 hours) leads us to suggest three mechanisms to help ensure that users will remember their action-object pairs in the Shared Cues password management scheme: 1) Start with a shorter initial time gap between the memorization phase and the first rehearsal (e.g., 3 hours or 6 hours). 2) Instruct the user to wait 12 hours after he has memorized the PAO stories before using the secret action-object pairs to form passwords. If the user can still remember his PAO stories after 12 hours then he can go ahead and use those stories to create passwords. 3) Implement a temporary recovery mechanism which allows a user who can remember one or two of his action-object pairs to recover his other action-object pairs during the first 24 hours (e.g., 98.6% of participants in the m_12hr×1.5 condition remembered their first action-object pair after 12 hours).

*Natural Rehearsals:* Following Blocki et al. [12] we provide a sense of extral rehearsal effort by assuming that the user's visitation schedule for each account is well-modeled by a Poisson arrival process, and that we know how frequently

| User | Number of Accounts Visited | | | $\mathbf{E}[XR_\infty]$ | $\mathbf{E}[XR_\infty - XR_{1.75}]$ |
|---|---|---|---|---|---|
| | Daily | Weekly | Monthly | | |
| Active | 5 | 5 | 4 | 3.29 | 0.01 |
| Typical | 2 | 8 | 4 | 7.81 | 0.14 |
| Infrequent | 0 | 2 | 12 | 30.41 | 7.41 |

TABLE IX: $\mathbf{E}[XR_\infty]$ — Expected number of extra rehearsals to remember 14 passwords with 4 PAO stories with $b = 0.5$ days and $\sigma_s = 1.5$

the user visits each of his accounts on average (e.g., daily, weekly, monthly). Table IX predicts how many extra rehearsals ($XR_\infty$) the user would need to do over his lifetime to ensure that he remembers all 4 of his PAO stories on average if the user initially returned after 12 hours and then in 1.5× increasing intervals (i.e., the schedule that provided the best result to the answer to the first question). The value $\mathbf{E}[XR_\infty]$ will depend on how frequently the user visits each of his accounts (e.g., $\mathbf{E}[XR_\infty]$ will be smaller for users who visits their accounts frequently because they are more likely to satisfy each of their rehearsal requirements naturally). The predictions indicate that more active users could maintain 14 secure passwords with minimal rehearsal effort. Table IX also predicts how many extra rehearsals the user would need to do after 1.75 days (the column labeled $\mathbf{E}[XR_\infty]$ - $\mathbf{E}[XR_{1.75}]$). We observe that most of the extra effort required of the user will come in the first few days — after 1.75 days we predict that our typical user will need to do 0.14 extra rehearsals over his lifetime to remember all four PAO stories.

*Mitigating the Interference Effect:* One potential downside of the Shared Cues password management scheme [12] is that the more secure versions of the scheme may require users to memorize multiple stories at once. For example, Blocki et al. [12] suggested that users memorize 43 stories to create 110 unique passwords with a $(43, 4, 1)$-sharing set family to ensure security against an offline adversary who has seen one or two of the user's passwords. In their scheme the user would need to memorize at least 36 of these stories to form the first 9 passwords. The results of our study indicate that it is more difficult for users to memorize more than two PAO stories at once (e.g., participants in the m_24hr×2+2start_1 and the m_24hr×2+2start_2 conditions were perfect while some participants in the m_24hr×2+2start_4 struggled to remember their action-object pairs). It is likely that the interference effect is, at least partially, due to user fatigue (e.g., participants who memorized four PAO stories had less mental energy to expend memorizing each action-object pair than participants who only memorized one or two PAO stories). One potential way to mitigate the interference effect would be have user's follow a staggered schedule in which they memorize two new PAO stories at a time. Another important research problem is to construct $(n, \ell, \gamma)$-sharing set families that expand gracefully so that the user does not need to memorize too many stories at the same time (e.g., for every $t$ we seek to memorize the number of action-object pairs that a user would need to memorize to form the first $t$ passwords).

REFERENCES

[1] J. Bonneau, C. Herley, P. C. van Oorschot, and F. Stajano, "The quest to replace passwords: A framework for comparative evaluation of web authentication schemes," in *IEEE Symposium on Security and Privacy*. IEEE, 2012, pp. 553–567.

[2] D. Florencio and C. Herley, "A large-scale study of web password habits," in *Proceedings of the 16th international conference on World Wide Web*. ACM, 2007, pp. 657–666.

[3] I. Center, "Consumer password worst practices," *Imperva (White Paper)*, 2010.

[4] H. Kruger, T. Steyn, B. Medlin, and L. Drevin, "An empirical assessment of factors impeding effective password management," *Journal of Information Privacy and Security*, vol. 4, no. 4, pp. 45–59, 2008.

[5] J. Bonneau, "The science of guessing: analyzing an anonymized corpus of 70 million passwords," in *Security and Privacy (SP), 2012 IEEE Symposium on*. IEEE, 2012, pp. 538–552.

[6] H. Ebbinghaus, *Memory: A contribution to experimental psychology*. Teachers college, Columbia university, 1913.

[7] P. Wozniak and E. J. Gorzelanczyk, "Optimization of repetition spacing in the practice of learning," *Acta neurobiologiae experimentalis*, vol. 54, pp. 59–59, 1994.

[8] P. Wozniak, "Supermemo 2004," *TESL EJ*, vol. 10, no. 4, 2007.

[9] P. Pimsleur, "A memory schedule," *The Modern Language Journal*, vol. 51, no. 2, pp. pp. 73–75, 1967. [Online]. Available: http://www.jstor.org/stable/321812

[10] A. Baddeley, *Human memory: Theory and practice*. Psychology Pr, 1997.

[11] L. STANDINGT, "Learning 10,000 pictures," *Quarterly Journal of Experimental Psychology*, vol. 5, no. 20, pp. 7–22, 1973.

[12] J. Blocki, M. Blum, and A. Datta, "Naturally rehearsing passwords," in *Advances in Cryptology - ASIACRYPT 2013*, ser. Lecture Notes in Computer Science, K. Sako and P. Sarkar, Eds. Springer Berlin Heidelberg, 2013, vol. 8270, pp. 361–380. [Online]. Available: http://dx.doi.org/10.1007/978-3-642-42045-0_19

[13] N. E. A. Guideline, "Electronic authentication guideline," April 2006.

[14] J. Bonneau and S. Schechter, ""toward reliable storage of 56-bit keys in human memory"," in *Proceedings of the 23rd USENIX Security Symposium*, August 2014.

[15] S. Designer, "John the Ripper," http://www.openwall.com/john/, 1996-2010.

[16] D. Goodin, "Why passwords have never been weaker-and crackers have never been stronger," http://arstechnica.com/security/2012/08/passwords-under-assault/, August 2012.

[17] "Zappos customer accounts breached," http://www.usatoday.com/tech/news/story/2012-01-16/mark-smith-zappos-breach-tips/52593484/1, January 2012, retrieved 5/22/2012.

[18] "Update on playstation network/qriocity services," http://blog.us.playstation.com/2011/04/22/update-on-playstation-network-qriocity-services/, April 2011, retrieved 5/22/2012.

[19] S. Biddle, "Anonymous leaks 90,000 military email accounts in latest antisec attack," http://gizmodo.com/5820049/anonymous-leaks-90000-military-email-accounts-in-latest-antisec-attack, July 2011, retrieved 8/16/2011.

[20] "An update on linkedin member passwords compromised," http://blog.linkedin.com/2012/06/06/linkedin-member-passwords-compromised/, June 2012, retrieved 9/27/2012.

[21] "Rockyou hack: From bad to worse," http://techcrunch.com/2009/12/14/rockyou-hack-security-myspace-facebook-passwords/, December 2009, retrieved 9/27/2012.

[22] "Data breach at ieee.org: 100k plaintext passwords," http://ieeelog.com/, September 2012, retrieved 9/27/2012.

[23] N. Provos and D. Mazieres, "Bcrypt algorithm."

[24] M. Fossi, E. Johnson, D. Turner, T. Mack, J. Blackbird, D. McKinney, M. K. Low, T. Adams, M. P. Laucht, and J. Gough, "Symantec report on the undergorund economy," November 2008, retrieved 1/8/2013. [Online]. Available: http://eval.symantec.com/mktginfo/enterprise/white_papers/b-whitepaper_underground_economy_report_11-2008-14525717.en-us.pdf

[25] J. Blocki, M. Blum, and A. Datta, "Naturally rehearsing passwords," *CoRR*, vol. abs/1302.5122, 2013.

[26] M. Buhrmester, T. Kwang, and S. D. Gosling, "Amazon's Mechanical Turk: A new source of inexpensive, yet high-quality, data?" *Persp. Psych. Sci.*, vol. 6, no. 1, pp. 3–5, 2011.

[27] P. G. Ipeirotis, "Demographics of Mechanical Turk," New York University, Tech. Rep. CeDER-10-01, 2010.

[28] J. S. Downs, M. B. Holbrook, S. Sheng, and L. F. Cranor, "Are your participants gaming the system? Screening Mechanical Turk workers," in *Proc. ACM CHI*, 2010.

[29] A. Kittur, E. H. Chi, and B. Suh, "Crowdsourcing user studies with Mechanical Turk," in *Proc. ACM CHI*, 2008.

[30] M. Toomim, T. Kriplean, C. Pörtner, and J. Landay, "Utility of human-computer interactions: toward a science of preference measurement," in *Proc. ACM CHI*, 2011.

[31] A. J. Berinsky, G. A. Huber, and G. S. Len, "Using Mechanical Turk as a subject recruitment tool for experimental research," *Political Analysis*, 2011.

[32] J. K. Goodman, C. E. Cryder, and A. A. Cheema, "Data collection in a flat world: Strengths and weaknesses of mechanical turk samples," *Journal of Behavioral Decision Making*, to appear.

[33] J. J. Horton, D. G. Rand, and R. J. Zeckhauser, "The online laboratory: Conducting experiments in a real labor market," *Experimental Economics*, 2010.

[34] E. Adar, "Why I hate Mechanical Turk research (and workshops)," in *Proc. CHI Workshop on Crowdsourcing and Human Computation*, 2011.

[35] R. S. Woodworth and H. Schlosberg, *Experimental psychology*. Oxford and IBH Publishing, 1954.

[36] "Claridge v. rockyou, inc." 2011.

[37] M. Burnett, *Perfect passwords: selection, protection, authentication*. Syngress Publishing, 2005.

[38] R. Monroe, "Xkcd: Password strength," http://www.xkcd.com/936/, retrieved 8/16/2011.

[39] S. Gaw and E. W. Felten, "Password management strategies for online accounts," in *Proceedings of the second symposium on Usable privacy and security*, ser. SOUPS '06. New York, NY, USA: ACM, 2006, pp. 44–55. [Online]. Available: http://doi.acm.org/10.1145/1143120.1143127

[40] J. Yan, A. Blackwell, R. Anderson, and A. Grant, "Password memorability and security: Empirical results," *Security & Privacy, IEEE*, vol. 2, no. 5, pp. 25–31, 2004.

[41] J. Stein, "Pimp my password," *Time*, p. 62, August 29 2011.

[42] S. Brand, "Department of defense password management guideline," 1985.

[43] K. Scarfone and M. Souppaya, "Guide to enterprise password management (draft)," *National Institute of Standards and Technology*, vol. 800-188, no. 6, p. 38, 2009.

[44] "Geek to live: Choose (and remember) great passwords," http://lifehacker.com/184773/geek-to-live--choose-and-remember-great-passwords, July 2006, retrieved 9/27/2012.

[45] D. Florêncio, C. Herley, and P. C. van Oorschot, "Password portfolios and the finite-effort user: Sustainably managing large numbers of accounts," in *Proceedings of the 23rd USENIX Security Symposium*, August 2014.

[46] R. Shay, P. Kelley, S. Komanduri, M. Mazurek, B. Ur, T. Vidas, L. Bauer, N. Christin, and L. Cranor, "Correct horse battery staple: Exploring the usability of system-assigned passphrases," in *Proceedings of the Eighth Symposium on Usable Privacy and Security*. ACM, 2012, p. 7.

[47] S. Komanduri, R. Shay, P. Kelley, M. Mazurek, L. Bauer, N. Christin, L. Cranor, and S. Egelman, "Of passwords and people: measuring the effect of password-composition policies," in *Proceedings of the 2011 annual conference on Human factors in computing systems*. ACM, 2011, pp. 2595–2604.

[48] J. Blocki, S. Komanduri, A. Procaccia, and O. Sheffet, "Optimizing password composition policies," in *Proceedings of the fourteenth ACM conference on Electronic commerce*. ACM, 2013, pp. 105–122.

[49] R. W. Proctor, M.-C. Lien, K.-P. L. Vu, E. E. Schultz, and G. Salvendy, "Improving computer security for authentication of users: Influence of proactive password restrictions," *Behavior Research Methods, Instruments, & Computers*, vol. 34, no. 2, pp. 163–169, 2002.

[50] Y. Zhang, F. Monrose, and M. K. Reiter, "The security of modern password expiration: an algorithmic framework and empirical analysis," in *Proceedings of the 17th ACM conference on Computer and communications security*. ACM, 2010, pp. 176–186.

[51] D. Florêncio and C. Herley, "Where do security policies come from,"
in *Proc. of SOUPS*, 2010, p. 10.

APPENDIX

*A. List of People, Actions and Objects from the User Study*

Here are a list of the people, actions and objects we used in the study. The lists contain 92 actions and 96 objects respectively.

**People:** Bill Gates, Bill Clinton, George W Bush, Lebron James, Kobe Bryant, Brad Pitt, Darth Vader, Luke Skywalker, Frodo, Gandalf, Michael Jordan, Tiger Woods, Michael Phelps, Angelina Jolie, Albert Einstein, Oprah Winfrey, Nelson Mandela, Bart Simpson, Homer Simpson, Adolf Hitler, Steve Jobs, Mark Zuckerberg, Justin Timberlake, Jay Z, Beyonce, Kim Jong Un, Joe Biden, Barack Obama, Pope Francis, Rand Paul, Ron Paul, Ben Afleck, Hillary Clinton, Jimmy Fallon

**Actions:** gnawing, mowing, rowing, oiling, egging, waving, bowing, seizing, stewing, signing, searing, bribing, swallowing, sucking, saving, sipping, tazing, tattooing, drying, dueling, dodging, tugging, taping, nosing, hunting, numbing, inhaling, knifing, nipping, muddying, miming, marrying, mauling, mashing, mugging, moving, mopping, racing, riding, reeling, reaching, raking, lassoing, welding, aligning, leashing, elbowing, juicing, shining, sheering, judging, choking, chipping, coating, concealing, destroying, kissing, aiming, kicking, punching, canning, combing, gluing, cooking, giving, copying, vising, voting, fanning, fuming, firing, fishing, high fiving, batting, burying, plowing, puking, popping, tasting, pulling, climbing, weeping, swimming, stretching, following, paddling, howling, smelling, rolling, waking, jumping

**Objects:** saw, teacup, hen, ammo, arrow, owl, shoe, cow, hoof, boa, sauce, suit, snow, piranha, chainsaw, shark, tiger, snake, razor-blade, sumo, seal, sock, safe, soap, daisy, toad, dime, tire, dish, duck, dove, ant, onion, wiener, nail, navy, menu, mummy, hammer, mail, microphone, horse, rat, iron, ram, pin, roach, rib, lion, lime, leach, lock, leaf, cheese, jet, chain, chime, gyro, chili, jeep, goose, cat, wagon, igloo, couch, cake, coffee, cab, vase, foot, phone, waffle, fish, bus, patty, bunny, bomb, pill, bush, bike, beehive, puppy, kite, canoe, boar, apple, moon, moose, tepee, ditch, key, shoe, home, toe, nose, cheetah