

Isn't that Fantabulous: Security, Linguistic and Usability Challenges of Pronounceable Tokens

Andrew M. White
amw@cs.unc.edu

Fabian Monrose
fabian@cs.unc.edu

Department of Computer Science

University of North Carolina at Chapel Hill

Katherine Shaw
kshaw@alumni.unc.edu

Elliott Moreton
moreton@email.unc.edu

Department of Linguistics

ABSTRACT

Over the past few decades, passwords as a means of user authentication have been consistently criticized by users and security analysts alike. However, password-based systems are ubiquitous and entrenched in modern society—users understand how to use them, system administrators are intimately familiar with their operation, and many robust frameworks exist to make deploying passwords simple. Unfortunately, much of the formal research on user authentication has focused on attempting to provide alternatives (e.g., biometrics) to password-based mechanisms (or belated analyses of users' password choices), forcing administrators to use ad-hoc methods in attempts to improve security. This practice has led to user frustration and inflated estimates of system security. We challenge common wisdom and re-examine whether *pronounceable* authentication strings might indeed offer a more reasonable alternative to traditional passwords. We argue that pronounceable authentication strings can lead to both improved system security and a decreased burden on users. To re-examine this potential, we explore questions related to how one might develop techniques for *rating* the pronounceability of word-like strings, and in doing so, enable one to quantify pronunciation difficulty. Armed with such an understanding, we posit new directions for *generating* usable passwords which are pronounceable and, we hope, memorable, hint-able and resistant to attack.

Categories and Subject Descriptors

K.6.5 [Management of Computing and Information Systems]: Security and Protection—*authentication*

General Terms

Security, Human Factors

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

NSPW'14, September 15–18, 2014, Victoria, BC, Canada.

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-3062-6/14/09 ...\$15.00.

<http://dx.doi.org/10.1145/2683467.2683470>

Keywords

pronounceable passwords; usable security; lexical blends

1. INTRODUCTION

Despite rampant criticism of passwords by users and security experts alike and an abundance of alternative proposals for user authentication, researchers have acknowledged that passwords are not likely to be replaced in the near future due to their ease of deployment and familiarity to users [23, 43]. Indeed, Bonneau et al. developed 25 different criteria intended to provide an objective viewpoint on the benefits and drawbacks of different user authentication schemes in terms of deployability, usability, and security; after evaluating passwords and 35 alternative schemes, Bonneau et al. ultimately concluded that no alternative scheme provides sufficient benefits to outweigh its detriments and overcome the current dominance of password-based systems [23].

We agree with that general sentiment, and believe it is time we refocus our attention on solutions that help improve the current state of affairs with password-based authentication. Over the past two decades, a number of policies for improving basic password systems have been suggested, but the most widely adopted of these suggestions is to increase the size of the space from which passwords are drawn (e.g., by enforcing the inclusion of numbers and special characters, requiring both upper and lower case letters, and increasing minimum password lengths). However, even for user-chosen secrets, these policies generally make passwords harder to remember and type, leading to user frustration and the habit of writing down and/or otherwise storing passwords [46, 92]. Worse yet, users generally fulfill these policy requirements in predictable ways, impairing, to a large extent, the security benefits these requirements are intended to provide [74, 87, 89]. Recently, Herley [42] rightfully argued that this behavior is, in fact, rational: users are often confronted with conflicting advice [54] and draconian password policies [46] for which they see marginal benefits [42]. Although alternatives to traditional textual passwords have been extensively explored and debated in the past [17]—particularly at NSPW [16, 37, 77, 80, 81]—we wish to stimulate further discussion and scientific exploration of methods which retain the benefits of textual passwords [23] while improving usability without sacrificing security.

In what follows, we argue that pronounceable authentication strings might offer a promising alternative to tra-

ditional passwords. In particular, we believe pronounceable authentication strings can lead to both improved system security and a decreased burden on users by providing memorable, hint-able passwords that are resistant to attack. We are intrigued by this direction because we believe early work on pronounceable passwords led to inconclusive results: while prior automated pronounceable password generators produced passwords which are both weak (due to their reliance on common English syllables) [55] and difficult to pronounce [28, 33], studies have also indicated that pronounceable passwords are easier for users to remember [75]. The latter finding suggests that the opportunity exists for leveraging linguistic expertise to develop techniques that provide memorable passwords through both user input and automated processes. One approach is to develop techniques for rating the pronounceability of word-like strings. In doing so, we believe we will be able to quantify pronunciation difficulty, which in turn, will allow us to proactively apply rigorous security analysis techniques to the space of pronounceable word-like strings to determine their suitability for use as secure passwords.

Our initial focus is on exploring the feasibility of automatically generating pronounceable passwords by forming lexical blends, sometimes known as *portmanteaus*, of two or more source words. By using pronounceable lexical blends, we believe these passwords will have a number of advantages over traditional systems, including *pronounceability*, *memorability*, *hintability*, and *resistance to attack*. We also envision scenarios where blends could be generated from one or more distinct semantic domains, which could be user-chosen. Generating blends from certain semantic domains naturally suggests one way in which we can provide *hintability*: the semantic domain(s) from which the source words are drawn can serve as hints (e.g., during the password reset process) to the user as to their password. Additionally, from a security perspective, system-generated lexical blends offer resistance to attacks as they will not be contained in any dictionary.

That said, we realize that in a proposition like ours, the security of the system depends on the space of possible blends; thus one must fully explore the security implications of that restriction. In particular, past research on analyzing the security of passwords has followed one of two approaches. The first is to consider only *system-generated* passwords, wherein the user is randomly assigned a password generated by the system in question. In this case, key properties to investigate include memorability and various negative externalities (such as password storing and user frustration) which result from passwords being assigned to users without any semantic context. The second approach is to consider empirically the choices of users, which we will refer to as *user-generated* passwords. Numerous other security-related questions abound, some of which we discuss herein.

Some of these questions might be addressed by an alternative type of password-generation, which we refer to as *user-influenced*, wherein a user can influence the formation of her password by, e.g., suggesting words to form part of a blend or suggesting semantic domains from which to draw source words, but the final password is ultimately system-generated. The hope is to aid in memorability by allowing the user to influence the formation of their (otherwise randomly generated) password. By allowing the system to ultimately assign the password, we may be able to maintain

enough uncertainty in the distribution of possible passwords to avoid compromising the security of the system.

Finally, we believe that the ability to generate and rate pronounceable strings has value that extends beyond the realm of authentication. For instance, the automatic generation of domain names which are pronounceable and consistently spellable would be a boon for marketing. Alternatively, generating pronounceable domain names that are similar to a given *seed* string might aid in the proactive registration of potentially typosquatted domains. In addition, the generation and rating of (consistently) pronounceable and consistently spellable strings from restricted semantic domains may enable, e.g., the names of new drugs to be semantically relevant and distinctive while remaining easy for consumers and health care practitioners alike to remember, spell, and pronounce. Other potential applications include the replacement of tracking numbers and transaction IDs with sets of distinct word-like strings which improve, e.g., recognition over the phone, while retaining the error-correcting qualities of today’s systems through metrics for measuring similarity in pronunciation. Finally, automatically generated yet pronounceable word-like strings might serve as an alternative to the ‘short authentication strings’ currently used to detect man-in-the-middle attacks on, e.g., VoIP calls following the ZRTP protocol [90].

2. RELATED WORK

Alternatives to memorizing passwords have been proposed, such as so-called *cognitive* [91] and *associative* [76] passwords. Cognitive passwords, also referred to as “personal knowledge questions”, are now often used as password-reset mechanisms. However, cognitive passwords suffer from some of the same weaknesses as traditional passwords: in particular, they are often easily circumvented by targeted attacks using information gathered through acquaintances, phishing, social media and other publically available data sources [25, 67, 70]. In addition, personal knowledge questions are subject to similar statistical attacks as traditional passwords [22].

Associative passwords are based on the premise that, when *cued* with a certain word (or list of words), a given individual will respond with the same associated word or words each time [76]. However, Bunnell et al. [25] suggested that associated words with low “guessability”—i.e., those known to produce a wide range of responses among different people—were actually *harder* for users to recall than randomly generated passwords. However, Bunnell et al. also emphasized that the appropriate experimental conditions for testing the usability of associative passwords had not yet been identified, and that further research was necessary to determine the factors influencing both the memorability and usability of associative passwords [25].

Another approach is the use of *passphrases*, consisting of multiple words rather than letters [66]. However, to date, there is little empirical evidence as to whether user-chosen passphrases offer any improvements over user-chosen passwords. Recent work suggests that users are unlikely to choose words with sufficient randomness to make passphrases secure against offline attacks [21]. In addition, a recent study by Shay et al., which compared system-assigned passphrases and passwords, found no significant benefit to either [75]. However, that study controlled for “guessability” by forcing relatively low entropies for the distributions from which passwords and passphrases were drawn, thus failing to take

advantage of the much larger space of words than characters. Whether the memorability of passphrases (of equal length) scales with the size of the space from which the words are drawn remains an open question.

3. PRONOUNCEABLE PASSWORDS

The notion of “pronounceable” passwords was first explored by Gasser [34] and later standardized by NIST in 1993. This scheme, which we refer to as APG, operates by sampling at random from a set of base units (e.g., individual letters and certain pairs of letters), combining the samples into valid syllables, then concatenating valid syllables to form a password. The probability distribution from which the base units are sampled is based on the unit frequencies in natural language; the concatenation at each step is governed by a complex set of rules also based on natural language [34].

Leonhard and Venkatakrisnan [55] proposed a pronounceable password generation scheme based on randomly selecting letters to form strings under two simple constraints (passwords may not begin or end with two consonants nor contain three consecutive consonants or vowels) intended to “ensure consistency with English spelling.” An advantage of their scheme is that the resulting password space is easy to analyze relative to schemes based on occurrence frequencies.

Other pronounceable password generators operate by building passwords using phonemes (distinct sounds), *digraphs* (i.e., pairs of letters representing distinct sounds) and/or *trigraphs* as building blocks, *n*-gram character models, and character (e.g., vowel-vowel) substitutions [28]. Unfortunately, the pronounceable password generators discussed in this section (including APG), have been criticized for producing passwords which are difficult to pronounce [28, 33]. These criticisms, however, have been subjective in nature. In Section 4, we present our preliminary work on an objective metric for measuring pronounceability, and apply this metric to passwords generated by five different ‘pronounceable’ password generators.

3.1 Open Security Considerations

When the distribution of (potential) passwords (or a large enough sample thereof) is available, statistical methods can be used to analyze the security of passwords randomly drawn from such a distribution. Early analyses of password security focused on an attacker model in which the attacker targeted a specific user. In recent years, however, it is clear that attackers are more interested in compromising as many accounts as they can rather than targeting specific users. Accordingly, security analyses have evolved from techniques which fail to accurately model password guessing difficulty in multi-account scenarios to more sophisticated metrics.

Dictionary Attacks. When passwords, or computable hashes thereof, are available, analysis of the passwords has often been performed using similar tools to those used by attackers. These tools generally incorporate a number of password dictionaries, often termed “wordlists”, and are capable of applying various transformations to the entries therein to provide variant passwords [54, 89]. The dictionaries used can be general or domain-specific: in particular, Kuo et al. explored the space of *mnemonic* passwords by collecting a dictionary of passwords formed by taking the first letter of each word in popular phrases (such as famous quotations or song lyrics). While their dictionary enabled them to crack a

smaller proportion of mnemonic passwords than a standard dictionary could ‘control’ passwords, their results nonetheless suggest that such dictionaries may reduce the effectiveness of mnemonic-based password schemes [54].

Therefore, it seems prudent to investigate the effectiveness of building dictionaries of pronounceable word-like strings and employ these dictionaries in simulated attacks. Although similar to the attacks performed previously on conventional passwords [89], the transformations applied to the various entries in the dictionary would include not only character transformations but also *phonetic* or *syllabic* transformations on the pronunciation of the password to mimic the password-generation process. As part of our preliminary work, we have also developed methods for generating lexical neighbors of existing pronunciations, which will provide the basis for these sorts of attacks. One such method is to substitute phonemes or syllables in pronounceable strings with randomly-chosen replacements, then check the resulting string for phonotactic correctness and pronounceability. We intend to investigate other methods of generating dictionaries for pronounceable word-like strings, as well.

Statistical Attacks. Ganesan and Davies proposed an attack on pronounceable password generators (specifically, a scheme implemented in Sandia’s Kerberos V distribution and APG) which exploits the difference between the probability of a password belonging to a particular ‘bucket’ (e.g., starting with a particular unit) and the proportion of all the possible passwords which belong to that bucket. By concentrating on passwords in small buckets with high likelihoods, the attacker gains an advantage. Although this attack is valid for the Sandia scheme, which employs a first step in which one of 25 templates is chosen uniformly at random but the distribution of the number of possible passwords across templates is non-uniform [33], the proposed bucketing for the APG scheme does not produce the same effect using the base unit distribution given by Gasser [34].

However, statistical attacks on pronounceable passwords have not been examined in great detail, leaving many directions open which would be prudent to explore. In particular, we intend to explore the value of applying α -*work-factor* [65] and α -*guesswork* [20] metrics to the distributions generated by pronounceable passwords schemes in order to quantify their security.

3.2 Open Usability Considerations

We believe that using pronounceable word-like strings as authentication tokens provides usability benefits over traditional password schemes, particularly in terms of reducing the cognitive load on users. Specifically, we hypothesize that, with respect to traditional passwords, pronounceable word-like strings:

1. are easier for users to remember (*memorability*)¹
2. are less frustrating for users (*acceptability*)
3. lead to fewer errors on token entry

Testing these hypotheses will require significant user studies for which rigorous and careful design is necessary. We

¹Interestingly, despite criticism on the grounds of pronounceability, Shay et al. [75] found that APG passwords were remembered slightly more often than passphrases.

believe the study of pronounceable passwords requires careful design decisions; many, but not all, of which have been explored and considered by at NSPW in the past (e.g., [16, 37, 80, 81]). Specifically, a number of important questions remain open to our minds, particularly when considered in the specific context of using pronounceable word-like strings for authentication.

1. Subjective evaluations of performance on memory-related tasks correlate poorly with objective evaluations in many domains [71]. We intend to verify whether this phenomenon persists in the specific domain of pronounceable word-like strings; if so, then the question becomes: *are subjective evaluations more important in promoting the adoption of a system as a whole than objective measurements?* Separate assessment of subjective and objective criteria ameliorates, to an extent, the bias introduced by this phenomenon. That said, the role of *perception* in such studies has been overlooked and deserves greater exploration. We hope to identify and bring to light other such factors that have been similarly disregarded in the past.
2. To what degree can studies in which participants are aware they are only using their password(s) for research provide ecological validity? As mentioned in Shay et al. [75], prior work has suggested that role-playing scenarios can influence users towards creating better passwords [47, 49]. But do better passwords imply ecological validity?
A recent study by Fahl et al. [31] compared passwords created by users for an online study to the real-world passwords employed by the users to protect their university accounts. The study found no significant difference, in similarity to participants’ real-world passwords, between primed (i.e., asked to simulate a high-value scenario) and non-primed (i.e., given no such suggestion) conditions [31]. Fahl et al.’s result suggests that such prompting may be unnecessary to achieve ecological validity; however, further studies are necessary before any strong conclusions can be drawn.
3. The results of any study comparing a new authentication method with traditional passwords are inherently biased by the relative unfamiliarity of users with the new method [17]. Any studies on pronounceable passwords and particularly lexical blends would have to overcome similar bias, though perhaps less so than, e.g., graphical passwords. However, while the bias toward the familiar may be small due to the similarity between pronounceable and traditional passwords, the same similarity also complicates the analysis by introducing ambiguity as to the actual extent of any bias.
4. Allowing users a significant period of time in which to become accustomed to the new system highlights the issue of ‘password interference’ [27], where the necessity of maintaining multiple passwords in memory contributes to degraded performance. The extent (and direction) of the bias this introduces to studies of pronounceable passwords—which are both similar to traditional passwords and yet distinct—is another important variable.

Taken as a whole, we see the first two issues above as open questions which deserve attention from the community. In addition, we believe the contribution of the two confounding variables cannot be ignored, despite likely biasing results

in opposing directions (a familiarity with textual passwords may induce a positive bias towards pronounceable passwords when compared with more radical alternatives while the effects of password interference may induce a negative bias). Since our proposed ideas differ in many ways from previous attempts to improve or replace traditional passwords, we solicit the community’s feedback on appropriate study design decisions in the context of using pronounceable word-like strings as authentication mechanisms.

4. PRELIMINARY IDEAS

For the purposes of this thought-exercise, we consider a *word-like* string as an ordered pair (s, p) consisting of a spelling and a pronunciation. For example, the word-like string *slig* consists of the spelling *slig* and the pronunciation $/slɪg/$ ². We consider a word-like string to be desirable based on three criteria, namely (1) **Wordlikeness**: It should look and sound normal for a word of that language, (2) **Consistent pronounceability**: Seeing the spelling s , all speakers should read it aloud in a similar way, as p , and (3) **Consistent spellability**: Hearing the pronunciation p , all speakers should infer the same spelling s . Wordlikeness, in terms of both lexical (text) and phonotactic (sound) similarity to known vocabulary, has been found to improve the memorability of a novel word [35, 36, 59, 60, 79]. Consistent pronounceability and spellability ensure that both the spoken and written forms of the string can be accurately learned from exposure to either of them. Examples of word-like strings satisfying or violating these criteria are shown in Table 1.

Rating Word-like Strings. Our current thinking for providing a rating metric is based on a simple scenario in which one person encounters an unfamiliar word in print and reads it aloud to another person, who then writes it down. The rating is a lower bound on the probability that what is written down by the listener is exactly what the speaker saw. In our preliminary proof-of-concept, we use a joint-sequence model of grapheme-to-phoneme and phoneme-to-grapheme conversion devised by Bisani and Ney [18]. A joint-sequence model is a model which maps sequences of symbols from one alphabet, e.g., *phonemes*, the distinct sounds which make up speech (such as the ‘b’ in ‘bat’), to sequences of symbols from another alphabet, e.g., *graphemes*, the distinct characters which form written language. The joint sequences in question are represented by a sequence of *graphones*, where each graphone is a sequence of zero or more letters paired with a sequence of zero or more phonetic symbols.

The training input to this model is simply a lexicon of words, i.e., (s, p) pairs. By forming all possible decompositions of each word-pronunciation pair into graphones, the model learns a probability distribution over graphone sequences using standard n -gram techniques. Given a novel written word s , we modified the model of Bisani and Ney to return a set of pronunciations $\{p_1, \dots, p_N\}$ together with associated conditional probabilities $\{\Pr(p_1 | s), \dots, \Pr(p_N | s)\}$. Due to its symmetry, the joint-sequence model functions equally well in the other direction, accepting a novel pronunciation p and returning spellings with associated probabili-

²In this paper, pronunciation is written using the International Phonetic Alphabet (IPA).

s	p	Wordlikeness	Consistent pronounceability	Consistent spellability
<i>snib</i>	/snɪb/	good	good	good
<i>fɪp</i>	/fɪp/	good	good	bad: <i>phɪp</i>
<i>smough</i>	/smoʊ/	good	bad: /smʌf/	bad: <i>smoe</i>
<i>tlib</i>	/tlib/	bad	good	good

Table 1: Examples illustrating the desirability criteria using word-like strings (of English). For those unfamiliar with IPA, /smoʊ/ rhymes with ‘dough’ while /smʌf/ rhymes with ‘tough’.

ties. In the following exposition, our model was trained on the Carnegie Mellon University Pronouncing Dictionary [86].

For our application, the desirability of a spelled candidate s is defined by its *maximum path probability*:

$$R(s) = \max_p \Pr(p | s) \Pr(s | p).$$

The path probability $\Pr(p | s) \Pr(s | p)$ is the probability that s will be read as p , and p written as s , as in the scenario above, and is therefore a lower bound on the probability that what the listener writes down is exactly what the speaker saw.

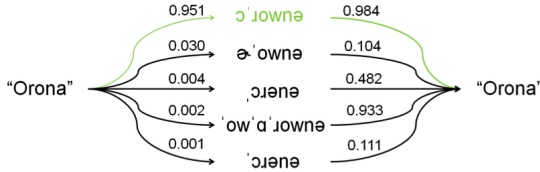


Figure 1: Possible paths for the string *Orona*.

Intuitively, graphemes which occur frequently in the training set are assigned higher probability by the model than those which are rare or absent. Thus, the model provides the means to recognize and distinguish English-looking and English-sounding novel words from others that are less English-like. The most direct way to exploit this capability would be to measure the wordlikeness of (s, p) as the probability assigned by the model to (s, p) . However, we used a simpler alternative, based on the observation that if a spelling s is atypical for English, there will be no conventional way to pronounce it and the model will be unable to parse the spelling into a sequence of graphemes which is clearly more probable than others. Consequently, none of the paths from s through some p back to s will have high probability, and so $R(s)$ will be low. The maximum path probability metric thus incorporates wordlikeness as well as consistency of pronunciation and spelling.

For a concrete example consider the candidate *orona*, some of whose paths appear in Figure 1. The numbers associated with the top path, highlighted in green, mean that *Orona* has a .951 probability of being pronounced as [ɔːˈrɔwnə], which in turn has a .984 probability of being written as *Orona*, according to the joint sequence model. Our preliminary rating metric provides us with the basic capability we need to rate the pronounceability of strings we generate. Obviously, this rating component can be used to filter the output that is shown to a user (in the case of system-generated authentication strings).

This preliminary metric also allows us to provide a quantitative analysis of previous work on pronounceable pass-

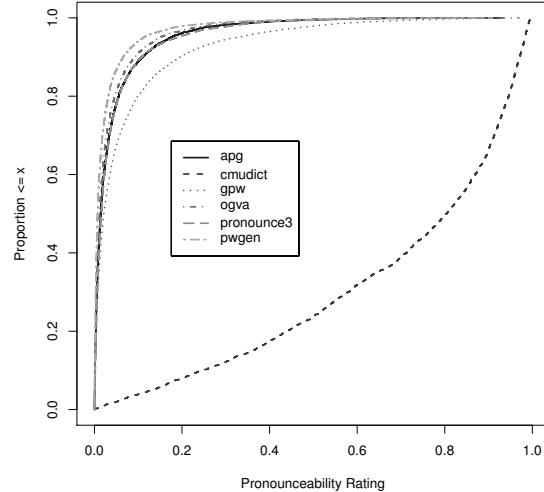


Figure 2: ECDF for ratings of 8-letter ‘pronounceable’ passwords from five different generators (apg, gpw, pwgen, pronounce3, and ogva; 5,000ea.) and 8-letter dictionary words (cmudict).

words. Our metric supports the previously mentioned criticism of ‘pronounceable’ password generators as producing difficult-to-pronounce passwords: as demonstrated in Figure 2, dictionary words (cmudict) score significantly higher using our metric than outputs from five different pronounceable password generators. The generators examined include two FIPS-181 implementations (apg³ and ogva⁴), a trigram-based algorithm (gpw⁵), a phoneme-based algorithm (pwgen⁶), and the character-based algorithm (pronounce3) proposed by Leonhard and Venkatakrisnan [55]. Perhaps unsurprisingly, the trigram-based gpw (which considers the relative frequencies in English of three-letter sequences representing distinct sounds) produces passwords which are, in aggregate, scored slightly higher than those of the other generators (which are based on smaller distinct units).

Estimating the Size of the Password Space. The pronounceable password generating algorithms mentioned above are necessarily complex due to the complicated set of rules

³<http://www.adel.nursat.kz/apg/> (v2.2.3)

⁴<https://pypi.python.org/pypi/passogva/1.0> (v1.0)

⁵<http://www.multicians.org/thvv/gpw.html>

⁶<http://pwgen.sourceforge.net/> (v2.05)

which governs spoken English. This complexity makes determining the resulting password spaces significantly more difficult than for traditional password schemes. An additional hurdle is that most of these algorithms are based on building passwords from component pieces that are spoken rather than written constructs. This necessitates a mapping from spoken form to written form before direct comparisons can be made.

To see why this mapping is important, consider a hypothetical generator which is based on stringing together n randomly selected syllables. The resulting written forms will vary in length significantly: in the CMU Pronouncing Dictionary [86], a one-syllable word’s written form can be up to 9 characters in length. A two-syllable word ranges up to 14 characters, and a three-syllable word up to 16 characters. Therefore, a direct comparison with, e.g., 8-character passwords is impossible.

The complexity of previous algorithms lead Leonhard and Venkatakrisnan [55] to design their algorithm with ease of analysis as a major goal. In addition, Gasser [34] performed simulations to determine the size of the password space for `apg` passwords of 6, 8, and 10 characters. The resulting values are given in Table 2 along with estimates of the size of each space if the generated passwords are filtered to retain only the top 10% in terms of pronounceability (see Figure 2).

Generator	Length	Original	Pronounceable
<code>apg</code>	8	2^{30}	2^{27}
<code>pronounce3</code>	8	2^{31}	2^{28}
<code>apg</code>	10	2^{40}	2^{36}

Table 2: Approximate password space sizes for pronounceable password generators with fixed password lengths (in characters). The final column reduces the size of the space to 10% of the original based on the results in Figure 2 and the assumption that marginally pronounceable passwords are filtered out.

Generating Word-like Strings. As part of our preliminary explorations, we have also developed a naïve generator for pronounceable word-like strings. The generator was designed to produce English-like output by building syllables from the sub-syllabic constituents (*onsets*, *nuclei*, and *rimes*) found in the words of the CMU Pronouncing Dictionary [86]. Candidates are generated by concatenating one to three syllables, each of which was made by concatenating a random onset, nucleus, and coda. Each candidate, being a pronunciation, is converted to its most-probable spelling using the trained joint-sequence model described above. This yields candidates with a variety of lengths and sound shapes.

Even with this naïve method, we can explore methods for producing lexical neighbors of real words, i.e., word-like strings which are close in pronunciation to existing words. For example, to generate new word-like strings that sound like a given “seed” word, we can take the seed word’s pronunciation and add, delete or replace phonemes to produce lexical neighbors. Candidates that are real words are removed from the set, as are those that cannot be parsed into legal syllables of English. The remaining candidates are converted to their most probable spellings and rated as before.

Table 3 shows a handful of top-rated candidates for several seed words.

ANIMATION	ENERGY	SOFTWARE
animationer	tenergy	siftware
anvimation	yenergy	sunftware
animationed	venergy	sulftware
ganimation	henergy	loftware
panimation	ebnergy	seftware

Table 3: Example seeded-candidates for some sample words.

Unfortunately, these seeded word-like strings would fail miserably as passwords for the same reason that traditional passwords should not resemble dictionary words, i.e., most cracking dictionaries and methodologies include variants of dictionary words as highly likely candidates for password guessing. They may also be somewhat difficult to remember. Therefore, we turn our attention to more advanced methods for producing word-like strings based on generating *lexical blends*. It is our hope that in doing so, we will arrive at a rich source of pronounceable word-like strings which are not found in dictionaries nor too ‘close’ to a single existing dictionary word.

In what follows, we explore a more ambitious possibility: generating nonwords which simultaneously sound like two given seed words, such as *thoughtomotive* (resembling *thought* and *automotive*), or *evergy* (resembling *energy* and *ever*). Alternatively, to avoid words that sound similar to those on a given “blacklist”, a filter stage can be added that compares the candidates to the blacklist and removes any that are closer than a user-specified distance.

4.1 Open Linguistic Considerations

Pronounceable password generators must be able to generate nonsense words that humans would nonetheless find appealing; however, the properties that make such word-like strings appealing also represent opportunities for an adversary to narrow down and possibly rank-order the space of possible passwords. Therefore, we are concerned with the linguistic factors that affect the security of pronounceable passwords. In the present work, our focus is on using *lexical blends* as passwords, and therefore we focus on the linguistic factors which affect lexical blend formation.

Previous research has identified several factors that influence blend formation, such as the phonemic content of the source words, their length in syllables, how they are stressed, their internal syllabic structure, and their frequency of use [e.g., 7, 13, 14, 26, 39, 48, 82–84]. What (to our surprise) has not been systematically studied is how blend formation is affected by *meaning*, i.e., by the semantic domain in which the blend is formed, or by the intended meaning of the blend—*information known by the password holder but less likely to be available to an adversary*. Therefore, an important prerequisite to systematic use of lexical blends as passwords is an understanding of how meaningful context can affect blend formation.

One of our primary objectives is to design and implement an automatic blend generator (henceforth “blender”), which produces blends similar to those produced by humans, to aid with this understanding. Simply stated, the goal is to correctly predict the pronunciation of a blend from those of

two given source words. In doing so, we hope to identify areas where there is variability in the pronunciation of a blend across different speakers. We believe it is important to focus on the variable blends because we hypothesize that two people blend the same word pair in different ways when they have different semantic definitions in mind. To simplify, we first consider candidates which begin like the first source word and end like the second (e.g., *brunch* from *breakfast* + *lunch*). Even so, there are several important questions that must be answered; we elaborate on some of them below.

Restricting the semantic domain. Blend formation is thought to be facilitated by *recoverability*, i.e., a blend is more acceptable to humans if the source words can be unambiguously inferred from it [13, 38]. For example, *education* + *entertainment* makes *edutainment*, not *educatement*, because *entertainment* is not recoverable from the latter — it could just as well be *payment*, *achievement*, etc. Recovery requires ruling out similar-sounding competitor words. Ordinary word recognition is faster and more accurate when context restricts the set of competitor words to associates (e.g., seeing *needle* facilitates recognition of *thread*) and relatives (e.g., *whale* and *dolphin*) of the target word [78]. We expect the same to be true of source-word recovery: blends will be more recoverable when they are processed in a restricted semantic domain. Hence, we expect meaningful context to make more blends possible, and to improve memorability.

We hypothesize that users can form more blends when helped by meaningful context. Our preliminary results test this hypothesis, which builds upon established methodologies in the ordinary word-recognition literature [e.g., 40, 61]. In particular, our pilot study explores several questions, including but not limited to: (1) does semantic context make source words more recoverable from blends? If so, then people should be faster and more accurate at “solving” blends when context situates them in a restricted semantic domain. For instance, will *tromboon* be solved faster (as *trombone* + *bassoon*) when the context is the related word *orchestra* than when it is the unrelated *alcohol* or when no context is given? Use of the same blend with different contexts controls for the inherent plausibility of the blended concept, as well as word frequency and phonological factors. (2) Are more blends possible in a restricted semantic domain? If recoverability facilitates blending, and a restricted search space facilitates recoverability, then we expect more blends to be possible when context establishes a semantically-restricted search space.

To see why these questions matter, consider a real-life example: The cast of characters in a book, movie, TV series, etc. is a restricted semantic domain, and fans refer to characters collectively by blending their names. For instance, *Minerva* McGonagall and *Severus* Snape, from the Harry Potter series, are collectively *Minerverus* or *Mineverus*. Such blends are often non-recoverable to an outsider, though obvious to an insider. Sometimes one character’s name is reduced to a single sound; e.g., *Kock* (*Kirk* plus *Spock*), *Rico* (*Rachel* plus *Nico*), or *Saylee* (*Simon* plus *Kaylee*) [85]. Blends that are opaque and non-viable without context can become viable when the context is known. Context may thus increase the number of possible blends, while also making them memorable for those who know the context and baffling for those who do not.

Structural ambiguity and ambi-blendability. Arndt-Lappe and Plag [7] have shown that many source-word pairs can be blended in more than one way. We hypothesize that the choice among blend options is affected by the intended meaning of the blend, as a function of its morphological and syntactic structure. In our preliminary analyses, we found several systematic ways of generating such *ambi-blendable* source pairs. One is to use words which match phonologically at two “pivots”, but diverge between them. For example, *flamingo* and *mongoose* match phonologically at the sound [m] and again at [ŋg]. The blend can switch words at either pivot, yielding two blend candidates — *flamongoose* and *flamingoose* — which preserve different amounts of the two source words. Other methods include unstressed medial syllables that agree in onset but disagree in rime, like *Hufflepuff* + *Gryffindor* → *Huffedor* or *Huffinpuff*, and onset clusters that could be switched at two points, like *blue* + *green* → *bleen* or *breen*.

We suspect that the user’s blending decision is affected by the intended meaning of the blend. Specifically, we hypothesize that the user’s morphosyntactic parse of the source-word pair affects the choice of blend via a phenomenon (called *positional privilege*) which has been extensively studied in the ordinary-language phonology of the world’s languages. Our recent work explores issues related to positional privilege in lexical blends [73]; our near-term goals include investigating how these issues might limit the space of pronounceable word-like strings that subjects create given two or more source words.

Estimating the Size of the Password Space. Similarly to simple pronounceable tokens, lexical blends present complications with respect to calculating the size of the password space. Previous work on lexical blends (e.g., [7, 13, 14, 26, 39, 48, 82–84]) has been from a bottom-up, inferential perspective, generally examining existing blends to determine the factors which result in the particular form that a blend of two source words takes and which is accepted into the language (e.g., why did blending ‘breakfast’ and ‘lunch’ result in ‘brunch’ rather than ‘breakfunch’?). While these factors, along with the intended meaning of the blend, are important to consider (particularly from a usability standpoint), they tell us little about the space of potential blends. What is needed, instead, is a top-down analysis of how often two source words can be combined in such a way as to produce a blend which is *acceptable*, in terms of phonotactic and syllabic constraints, regardless of whether it’s the *preferred* blend. Unfortunately, the complexity of spoken English makes this type of analysis difficult without extensive simulations.

To see why extensive simulations are necessary, it is instructive to consider the many different forms which blends can take [72]. For instance, the combination of two source words can leave both words intact (as in *exam* + *amnesia* → *examnesia*), leave one word intact but use only part of the other word (as in *decathlon* + *athlete* → *decathlete*), or leave neither word intact (as in *motor* + *hotel* → *motel*). Furthermore, the two source words can share segments, i.e., overlap, in the blend (as in the examples in the previous sentence) or share nothing (as in *breakfast* + *lunch* → *brunch*). Finally, these examples cover only *linear* blend structures: other structures, such as embedded blends like *advertisement* + *tease* → *advert~~te~~asement*, exist. Each of these pos-

sibilities must be investigated for each pair of source words and, for every possible pivot point or shared segment, the result tested against the phonotactic and syllabic constraints of spoken English.

The previous paragraph notwithstanding, we can make some rough estimates from our preliminary work on an automatic blend generator (our “blender”) [72]. In order to ease the linguistic analyses, our preliminary blender is narrowly focused on one specific type of blend created from a particular subset of available word-pairs. As such, the result which follows is strictly a lower-bound on the number of possible blends and is likely to be extremely pessimistic. In particular, we used as our source of words the intersection of the CMUdict (as syllabified by Bartlett et al. [11]) and CELEX [9] lexicons (in order to eliminate many of the proper nouns in CMUdict). This resulted in a list of 36,216 words. We further restricted our attention to nouns consisting of exactly two syllables. From the possible pairs of words remaining, we selected those with a shared consonant between the nuclei of the two syllables to use as the blending point. This ensures that the blend would be phonotactically viable, resulting in 2,600,363 blends. We remind the reader that this figure is based only on pairs of disyllabic nouns blended at a shared internal consonant and is therefore an exceedingly loose lower-bound on the number of possible blends. That said, forthcoming work [32] suggests that 10^6 guesses provide a reasonable limit on the capabilities of an online attacker; in expectation, even the lower-bound given above surpasses this limit (assuming the blends are chosen uniformly at random). Security against online attack, therefore, appears to be an attainable goal in the system-generated model.

Blends in Languages other than English. Our analysis and discussion so far has focused on English. The formation and usage of blends, however, differs from one language to another. Table 4 shows the 15 languages with the most native speakers [57]. Their combined speaker population is 3.6 billion, which is more than half of the world’s population. Seven of these languages—Spanish, English, Portuguese, Japanese, German, Korean, and French, with a total of 1.3 billion speakers—are known to have productive lexical blending, i.e., existing blends are abundant, and speakers readily invent new ones. In four others—Chinese, Hindi, Arabic, and Russian, with 1.9 billion speakers—blends are rare, though not necessarily non-existent [69]. We found no data for the four remaining languages: Bengali, Lahnda (Punjabi and its relatives), Javanese, and Telugu.

The fact that blending is not productive in a language does not necessarily mean that speakers cannot readily adopt the process. In recent years, blending has spread to languages which previously had few or no blends. For instance, until the late 20th century, Polish (a close relative of Russian, with 39 million native speakers) had only a handful of lexical blends, which were not recognized as such by most speakers. However, since the 1960s, and especially the 1990s, blending has become a major source of new Polish words under the influence of English and other international languages [51]. The same has happened in Korean (77 million) [2], Ukrainian (32 million) [24], Hebrew (4.8 million) [19], and Greek (11 million) [68]. Speakers of Hindi can judge the well-formedness of experimenter-devised Hindi blends, even though naturally-occurring blends are rare or

nonexistent in their language [63]. It may therefore be possible to use blending as a password scheme even in languages which do not currently have productive blending.

In all languages which do have productive blending, there are multiple blend formation strategies (overlapping, insertion of one word into another, splicing at a syllable boundary, etc.). The choice of strategy for any particular source-word pair is determined by multiple interacting phonological, morphological, and semantic factors, as well as by avoidance of collision between the new blend and existing words [8, 13, 39]. A separate question is whether English is unique in having source-word pairs that can be blended in more than one way (*ambi-blendability*), like *blue* + *green* → *bleen* or *grue*. We know of no study that focuses on this question, but examples have been reported in languages as far apart as Hebrew [12] and Korean [2].

4.2 Possible Extensions & Associated Challenges

User Influence on Password Formation. Allowing user input into the password formation process fundamentally changes the way in which we must analyze password-based systems. In order to study these changes, we propose the use of both small, lab-based studies and larger scale crowd-sourced studies to generate pronounceable passwords under a variety of user-influenced conditions, as mentioned earlier. The passwords generated in these cases will then be subject to the same rigorous analyses as in the case of (fully) system-generated passwords. In addition, we fully anticipate that new avenues of attack will be unveiled during the course of this project, such as trends in user inputs which might provide an advantage to an attacker.

Prior work has analyzed biases in user-chosen passwords in both traditional and graphical password schemes [29, 89]. The work in Zhang et al. [89], for example, identified a number of distinct transforms which were common across many users, as well as identifying and exploiting the habits of particular users. For that reason, we believe it is also important to investigate whether our user-influenced password generation techniques lead users to form similar habits in password formation and/or generate passwords which are predictable from past passwords. Our prior work on ‘phonetic edit distance’ metrics [88] positions us perfectly to perform this analysis on pronounceable passwords by allowing us to measure the similarity of two passwords based on phonetic transformations as well as character transformations.

Unfortunately, even relatively naïve schemes for reducing the size of the potential pronounceable password space, such as rejecting any potential passwords which violate basic linguistic constraints, may substantially improve an attacker’s chances. Therefore, a thorough analysis of pronounceable passwords and lexical blends is necessary and must include worst-case estimates. For instance, we intend to explore whether there exist hitherto unknown paradigms, e.g., dominant strategies in user formation of lexical blends, that an attacker might exploit. We also intend to assess the extent to which the use of pronounceable passwords and lexical blends impacts the effectiveness of shoulder-surfing. We welcome suggestions for other such dimensions to investigate.

Second-order Hints as a Password Reset Mechanism. While we suspect that the number of consistent blends for a given pair of source words will be too small to facilitate

Language	Speakers	Blends	Example	References
Chinese	1197M	No	<i>dīng kè zú</i> ‘DINKs’ + <i>gǒu</i> ‘dog’ → <i>dīnggǒuzú</i> ‘DINKs with a dog’	Ronneberger-Sibold [69]
Spanish	414M	Yes	<i>loca</i> ‘crazy’ + <i>Colombia</i> → <i>Lo<u>co</u>mbia</i> ‘crazy Columbia’	Piñeros [64]
English	335M	Yes	<i>spoon</i> + <i>fork</i> → <i>spork</i>	Algeo [5]
Hindi	260M	No	—	Ohala [63]
Arabic	237M	No	—	Al-Hamly and Farghal [4]
Portuguese	203M	Yes	<i>telefone</i> + <i>móvel</i> ‘mobile’ → <i>tele<u>mó</u>vel</i> ‘cell phone’	de Araújo [30]
Bengali	193M	—	—	—
Russian	167M	No	—	Arcodia and Montermini [6]
Japanese	122M	Yes	<i>gorira</i> ‘gorilla’ + <i>kuzira</i> ‘whale’ → <i>Goz<u>ira</u></i> ‘Godzilla’	Kubozono [53]
Javanese	84M	—	—	—
Lahnda	82M	—	—	—
German	78M	Yes	<i>Tomate</i> ‘tomato’ + <i>Kartoffel</i> ‘potato’ → <i>Tom<u>off</u>el</i> ‘so-matic hybrid of tomato and potato’	Ronneberger-Sibold [69]
Korean	77M	Yes	<i>p^hok^hi</i> ‘fork’ + <i>sutkalak</i> ‘spoon’ → <i>p^hok^hal<u>ak</u></i> ‘spork’	Ahn [3]
French	75M	Yes	<i>chien</i> ‘dog’ + <i>chimpanzé</i> ‘chimpanzee’ → <i>ch<u>ien</u>-pan<u>zé</u></i> ‘a dog that acts like a chimpanzee’	Bertinetto [15], Léturgie [56]
Telugu	74M	—	—	—

Table 4: Blend productivity in languages with the most native speakers. Population statistics [1] for Chinese include mutually-unintelligible non-Mandarin languages; Lahnda includes Punjabi and related languages.

password-reset (or hinting) mechanisms if we simply show the user the two words and ask her to provide the “correct” blend, there are other possible ways to provide hints to a password. One method we intend to explore is to provide users with one of the source words. Others include (1) remind the user of the semantic domains of the source words, and (2) provide the user with hints whose connection to the source word is obscure to anyone who does not know the source word semantic domains. These methods are similar to *associative* passwords.

Consider, for example, the scenario where a user is asked for areas of interests during password setup. For pedagogical purposes, assume she specifies the semantic domains *genetics* and *knitting*. As output, the system-generated pronounceable password “*homeosocks*” (from *homeobox* plus *socks*) is suggested, which she later accepts as her password. In the event that she forgets her password, the system could prompt her with *genetics* or *knitting*, which could jog her memory without giving the word away to an attacker. Less directly, it can prompt her with *regulate* or *feet*, which are associates of the target words only to someone who is already thinking in those domains. The user is likely to interpret the hints in the correct semantic domains (which she chose herself), allowing her to retrieve the hinted source word and hence the blend. For an adversary, the hints point in many more directions (e.g., *regulator* could lead to a fruitless search among scuba-diving terms).

On the other hand, ambi-blendability could result in blend passwords which are harder to hint, since the user might recall both source words, but still choose the wrong blend for them. Additionally, we found that 469,190 ambi-blendable word pairs of one very specific type could be made from a dictionary of 36,216 words, suggesting that ambi-blendable pairs (to say nothing of triplets or longer word strings) may be plentiful. Our results, however, also suggest precautions against this danger: by exploiting positional-privilege effects, hints could steer the user in the right direction with semantic context. For example, *Has a tail, steals things*

is an appropriate hint for *baboondit*, whereas *Steals things with a tail* is appropriate for *babandit*—and not the other way around. Even an unsophisticated hint generator, such as that discussed previously, could make use of simple hint schemas (e.g., “X and Y” vs. “X of Y”) to nudge the user one way or the other.

Natural Second Factors. There are at least two natural second factors for pronounceable-password-based authentication mechanisms. Since these passwords are inherently pronounceable, voice-based authentication suggests itself as a likely candidate. In addition, verbal input may be a viable alternative input mechanism (e.g., for mobile devices or over-the-phone exchanges) where text input is unavailable or inconvenient. The second natural candidate for two-factor authentication is *lip movement* as a biometric, especially because we can ensure that our system-generated passwords induce a threshold number of lip movements. Finally, prior work suggests that the combination of voice analysis and lip movement provides a reasonable biometric for speaker recognition [45]. Further analysis is necessary to determine which of these second factors might be the most fruitful to explore in the short term.

4.3 Preliminary Experiment Designs

Testing the hypotheses present in this work will require substantial experimentation, particularly in the form of user studies. In this section, we outline our thoughts on the major security-related experiments to perform.

System-Assigned Tokens. Our first experiment would concern system-assigned authentication tokens, including traditional passwords, pass-phrases, pronounceable passwords, and lexical blends. We intend to build upon the work of Shay et al. [75], who performed a laudably thorough experiment and subsequent analysis, and their obvious and commend-

able commitment to full disclosure and reproducibility will no doubt aid us in our own work.

There are, however, aspects of Shay et al.’s experimental design which we disagree with and which we intend to re-think. In particular, Shay et al. used a fixed entropy value (of 30 bits) to determine the spaces from which the various authentication tokens they used were drawn. While this allows for straightforward comparison in terms of the security provided by these schemes, it mitigates the advantages of schemes which have a naturally higher entropy for the same cognitive requirements on the part of the user. For instance, from the point of view of a user, remembering a phrase consisting of three words each chosen at random from a dictionary of 200 words is no less difficult than remembering a similar phrase with words chosen from a dictionary of 10,000 words (provided the word list is not available to the user). Shay et al.’s study supports this contention in that the size of the dictionary (181, 401, or 1,024 words in their study) resulted in no significant difference in successful recall rate.

We suggest, therefore, that the different authentication schemes tested should not be artificially limited to a specific level of security, but rather allowed to assume the level of security provided in their ‘natural’ configuration. For instance, the median adult vocabulary size varies by age between 25,000 and 30,000 words [1]. Pass-phrase schemes with two, three, or four words chosen uniformly at random from a dictionary of size 25,000 have entropies of 29.2, 43.8, and 58.4 bits, respectively, the latter two a significant improvement over the 30 bits used as a baseline by Shay et al. without, we argue, a significant increase in cognitive load on the part of the user.

One might now argue that comparisons between schemes are no longer fair, since the security level of the system is no longer fixed. One approach to mitigating this problem is to compare each non-password scheme against a password scheme with the same level of security. Another approach would be to make relatively minor adjustments to the various schemes so as provide a small number of groups with roughly equivalent security. Forthcoming work [32] suggests that only two such groups are necessary: schemes providing tokens which will, in expectation, survive 10^6 guesses (sufficient to withstand online attacks) and 10^{14} guesses (for offline attacks). Finally, statistical measures, such as ANCOVA (analysis of covariance) [44], can control for the difference in security.

In addition to traditional passwords and pass-phrases, we intend to include system-assigned pronounceable tokens (and request pronounceability measurements from the users), lexical blends, and lexical blends of words chosen from restricted semantic domains. For the latter, conditions might include notifying the user of the semantic domain, using the semantic domain as a hint, and withholding the domain from the user. A comparison with pass-phrases under similar selection and presentation conditions is also warranted.

User-Influenced Tokens. Our second major experiment would focus on user-influenced tokens. Treatments for this experiment might include user-created lexical blends under various conditions: without any prompting (beyond an explanation of the general scheme), with a given semantic domain (possibly chosen by the user), with one word assigned by the system, and with the user choosing from a list of

system-generated blends. For the sake of comparison, pass-phrase treatments might be included under similar conditions. Further conditions might include various second-order hint mechanisms, such as providing the user with a hint like the semantic domain, another word from the same semantic domain, or one of the source words. One primary goal would be to assess users’ relative ability to recall tokens across different schemes. Additionally, users’ preferences would be assessed. Finally, another primary goal would be to elicit an empirical distribution for each type of authentication token, to be used in analyses of the effective security levels such as have previously been performed for passwords (e.g., [20, 87]) and pass-phrases [21]. The results of these analyses could be used to inform any recommendations based on recall rates.

5. CONCLUSION

With this work, we hope to call attention to both the promise and the challenges inherent in exploring the use of pronounceable word-like strings, and particularly lexical blends, as passwords. We outlined approaches to both rating the pronounceability of word-like strings and their automated generation, which we argue are essential tools for scientifically analyzing the feasibility of pronounceable passwords. Also central to that effort is an understanding of the linguistic properties of pronounceable passwords. Towards that end, we reported on our preliminary investigations into lexical blends, which we believe may serve as a viable pool of pronounceable, memorable and hintable passwords which are resistant to attack. Finally, we solicited the community’s feedback on designing appropriate experiments to test our hypotheses, and highlighted a number of potential issues in experimental design which we believe are particularly relevant in the context of pronounceable passwords.

The questions and challenges discussed in the paper are by no mean the only issues worth exploring with regards to the applicability of user-influenced, yet system generated, pronounceable tokens. For instance, it is conceivable that allowing for error correction on the generated tokens might provide a substantial usability benefit. Indeed, this was already shown to be the case for passphrases, where users often misspell words when attempting to recall their passphrases [75]. In the case of pronounceable passwords, such error correction can be performed not only at the orthographic level but at the pronunciation level as well; e.g., ‘piece’ and ‘peace’ are pronounced identically but spelled differently. This sort of error correction may lead to increased usability for pronounceable passwords.

That said, incorporating error correction is a double-edged sword, as attackers may gain an advantage from the practice. Therefore, the limits of this advantage, from both phonetic and orthographic perspectives, must also be explored. We suspect that it may be possible to mitigate the attacker’s advantage by incorporating password distinctiveness (phonetically, orthographically, or both) into the generation process (be it system generated or user-influenced). We may be able to use phonetic edit distance metrics [88] or lexical neighborhoods to determine whether isolating passwords in this way is feasible and to what extent this will reduce the space of pronounceable passwords.

It is also conceivable that some applications will call for novel words which resemble specific real words, while others may demand that the novel words *avoid* sounding like particular real words. Both desiderata require a metric of word

similarity. One idea is to use the Levenshtein (string-edit) distance between pronunciations, which is defined as the minimum number of phoneme insertions, deletions, and substitutions required to turn one pronunciation into another [52]. In particular, the string-edit distance has been shown to influence perceived word similarity in psycholinguistic experiments [10, 41, 58] and in real-world confusions [50]. However, this glosses over the fact that edit distance does better if the different elementary operations are weighted differently, so that substituting similar-sounding phonemes contributes less distance than if the phonemes were distinct. Nevertheless, this issue is also worth exploring.

6. ACKNOWLEDGMENTS

We are thankful to Steve Bellovin, Matt Bishop, Joseph Bonneau, and Paul van Oorschot for insightful discussions and comments on earlier drafts of this manuscript. We would also like to thank the NSPW 2014 participants for their feedback and Bob Blakley, in particular, for his handwritten transcription of the discussion. This work is supported in part by a grant from the National Science Foundation, under award number 1318520.

7. REFERENCES

- [1] Native speakers in greater detail, May 2013. Retrieved July 30, 2014 from <http://testyourvocab.com/blog/2013-05-08-Native-speakers-in-greater-detail>.
- [2] S. Ahn. A constraint-based analysis of Korean blends. Master's thesis, Seoul National University, 2012.
- [3] S. Ahn. Faithfulness conflict in Korean blends. *University of Pennsylvania Working Papers in Linguistics*, 20(1):1–10, 2013.
- [4] M. Al-Hamly and M. Farghal. English reduced forms in Arabic scientific translation: a case study. *Jordan Journal of Modern Languages and Literature*, 5(1): 1–18, 2013.
- [5] J. Algeo. Blends, a structural and systemic view. *American Speech*, 52(1/2):47–64, 1977.
- [6] G. F. Arcodia and F. Montermini. Are reduced compounds compounds? Morphological and prosodic properties of reduced compounds in Russian and Mandarin Chinese. In V. Renner, F. Maniez, and P. J. L. Arnaud, editors, *Cross-disciplinary perspectives on lexical blending*, pages 94–113. de Gruyter Mouton, Berlin, 2012.
- [7] S. Arndt-Lappe and I. Plag. Phonological variability in English blends. Handout from the Workshop on Data-Rich Approaches to English Morphology, Victoria University, Wellington, New Zealand, July 4–6 2012.
- [8] S. Arndt-Lappe and I. Plag. The role of prosodic structure in the formation of English blends. *English Language and Linguistics*, 17(3):537–563, 2013.
- [9] R. Baayen, R. Piepenbrock, and L. Gulikers. CELEX2 LDC96L14. Web Download. Philadelphia: Linguistic Data Consortium, 1995.
- [10] T. Bailey and U. Hahn. Determinants of wordlikeness: Phonotactics or lexical neighborhoods. *Journal of Memory and Language*, 44(4):568–591, May 2001.
- [11] S. Bartlett, G. Kondrak, and C. Cherry. On the syllabification of phonemes. In *Proceedings of Human Language Technologies*, 2009.
- [12] O. Bat-El. Selecting the best of the worst: the grammar of Hebrew blends. *Phonology*, 13:283–328, 1996.
- [13] O. Bat-El. Blend. In K. Brown, editor, *The Encyclopedia of Language and Linguistics*, volume 2, pages 66–70. Elsevier, Oxford, England, 2nd edition, 2006.
- [14] O. Bat-El and E.-G. Cohen. Stress in English blends: a constraint-based approach. In V. Renner, F. Maniez, and P. J. L. Arnaud, editors, *Cross-disciplinary perspectives on lexical blending*, pages 193–211. de Gruyter Mouton, Berlin, 2012.
- [15] P. M. Bertinetto. Blends and syllable structure: a four-fold comparison. In M. Lorente, N. Alturo, E. Boix, M. R. Loret, and L. Payrató, editors, *La gramàtica i la semàntica per a l'estudi de la variació*. PPU-Secció de Lingüística Catalana de la Universitat de Barcelona, Barcelona, 2001.
- [16] K. Bicakci and P. C. van Oorschot. A multi-word password proposal (gridword) and exploring questions about science in security research and usable security evaluation. In *Workshop on New Security Paradigms*, 2011.
- [17] R. Biddle, S. Chiasson, and P. Van Oorschot. Graphical passwords: Learning from the first twelve years. *ACM Computing Surveys*, 44(4):19:1–19:41, Sept. 2012.
- [18] M. Bisani and H. Ney. Joint-sequence models for grapheme-to-phoneme conversion. *Speech Communication*, 50(5):434–451, May 2008.
- [19] S. Bolozky. *Measuring productivity in word formation: The case of Israeli Hebrew*. Brill, 1999.
- [20] J. Bonneau. The science of guessing: analyzing an anonymized corpus of 70 million passwords. In *IEEE Symposium on Security & Privacy*, 2012.
- [21] J. Bonneau and E. Shutova. Linguistic properties of multi-word passphrases. In *Workshop on Usable Security*, 2012.
- [22] J. Bonneau, M. Just, and G. Matthews. What's in a name? In *Financial Cryptography and Data Security*, volume 6052 of *Lecture Notes in Computer Science*, pages 98–113. Springer, 2010.
- [23] 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 & Privacy*, 2012.
- [24] S. R. Borgwaldt, T. Kulish, and A. Bose. Ukrainian blends: elicitation paradigm and structural analysis. In V. Renner, F. Maniez, and P. J. L. Arnaud, editors, *Cross-disciplinary perspectives on lexical blending*, pages 75–92. de Gruyter Mouton, Berlin, 2012.
- [25] J. Bunnell, J. Podd, R. Henderson, R. Napier, and J. Kennedy-Moffat. Cognitive, associative and conventional passwords: Recall and guessing rates. *Computers & Security*, 16(7):629–641, 1997.
- [26] G. Cannon. Blends in English word formation. *Linguistics*, 24:725–753, 1986.
- [27] S. Chiasson, A. Forget, E. Stobert, P. C. van Oorschot, and R. Biddle. Multiple password interference in text passwords and click-based graphical passwords. In *ACM Conference on Computer and Communications Security*, 2009.

- [28] H. Crawford and J. Aycock. Kwyjibo: automatic domain name generation. *Software: Practice and Experience*, 2008.
- [29] D. Davis, F. Monrose, and M. K. Reiter. On user choice in graphical passwords schemes. In *USENIX Security Symposium*, 2004.
- [30] G. A. de Araújo. Morfologia não-concatenativa em português: os portmanteaux. *Cadernos de Estudos Lingüísticos*, 39:5–21, 2000.
- [31] S. Fahl, M. Harbach, Y. Acar, and M. Smith. On the ecological validity of a password study. In *Symposium on Usable Security and Privacy*, 2013.
- [32] D. Florêncio, C. Herley, and P. C. van Oorschot. An administrator’s guide to internet password research. In *Large Installation System Administration Conference*. USENIX Association, Nov. 2014. Forthcoming.
- [33] R. Ganesan and C. Davies. A new attack on random pronounceable password generators. In *NIST National Computer Security Conference (NCSC)*, 1994.
- [34] M. Gasser. A random word generator for pronounceable passwords. Technical Report MTR-3006, MITRE Corporation, 1975.
- [35] S. E. Gathercole, C. R. Frankish, S. J. Picering, and S. Peaker. Phonotactic influences on short-term memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 25(1):84–95, 1999.
- [36] S. E. Gathercole, C. R. Frankish, S. J. Picering, and S. Peaker. Correction to Gathercole et al. (1999). *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 25(3):562, 1999.
- [37] M. Gibson, K. Renaud, M. Conrad, and C. Maple. Musipass: authenticating me softly with “my” song. In *Workshop on New Security Paradigms*, 2009.
- [38] S. T. Gries. Isn’t that *fantabulous*? How similarity motivates intentional morphological blends in English. In M. Acard and S. Kemmer, editors, *Language, Culture, and Mind*, chapter 28, pages 415–428. CSLI Publications, Stanford, California, 2004.
- [39] S. T. Gries. Shouldn’t it be *breakfunch*? A quantitative analysis of blend structure in English. *Linguistics*, 42(3):639–667, 2004.
- [40] F. Grosjean and U. Frauenfelder. A guide to spoken word recognition paradigms: Introduction. *Language and Cognitive Processes*, 11(6):553–558, 1996.
- [41] U. Hahn and T. M. Bailey. What makes words sound similar? *Cognition*, 97:227–267, 2005.
- [42] C. Herley. So long, and no thanks for the externalities: The rational rejection of security advice by users. In *Workshop on New Security Paradigms*, 2009.
- [43] C. Herley and P. Van Oorschot. A research agenda acknowledging the persistence of passwords. *IEEE Security & Privacy Magazine*, 10(1):28–36, 2012.
- [44] B. E. Huitema. Analysis of covariance (ANCOVA). In K. R. Neil J. Salkind, editor, *Encyclopedia of Measurement and Statistics*. Sage Publications, Inc., 2007.
- [45] M. Ichino, H. Sakano, and N. Komatsu. Multimodal biometrics of lip movements and voice using kernel fisher discriminant analysis. In *International Conference on Control, Automation, Robotics and Vision*, 2006.
- [46] P. G. Inglesant and M. A. Sasse. The true cost of unusable password policies. In *SIGCHI Conference on Human Factors in Computing Systems*, 2010.
- [47] P. G. Kelley, S. Komanduri, M. L. Mazurek, R. Shay, T. Vidas, L. Bauer, N. Christin, L. F. Cranor, and J. Lopez. Guess again (and again and again): Measuring password strength by simulating password-cracking algorithms. In *IEEE Symposium on Security & Privacy*, 2012.
- [48] M. H. Kelly. To “brunch” or to “brench”: some aspects of blend structure. *Linguistics*, 36(3):579–590, 1998.
- [49] S. Komanduri, R. Shay, P. G. Kelley, M. L. Mazurek, L. Bauer, N. Christin, L. F. Cranor, and S. Egelman. Of passwords and people: measuring the effect of password-composition policies. In *SIGCHI Conference on Human Factors in Computing Systems*, 2011.
- [50] G. Kondrak and B. Dorr. Automatic identification of confusable drug names. *Artificial Intelligence in Medicine*, 36(1):29–42, 2006.
- [51] E. Konieczna. Lexical blending in Polish: a result of the internationalisation of Slavic languages. In V. Renner, F. Maniez, and P. J. L. Arnaud, editors, *Cross-disciplinary perspectives on lexical blending*, pages 51–73. de Gruyter Mouton, Berlin, 2012.
- [52] J. B. Kruskal. An overview of sequence comparison: time warps, string edits, and macromolecules. *SIAM Review*, 25(2):201–237, 1983.
- [53] H. Kubozono. The mora and syllable structure in Japanese: evidence from speech errors. *Language and Speech*, 32(3):249–278, 1989.
- [54] C. Kuo, S. Romanosky, and L. F. Cranor. Human selection of mnemonic phrase-based passwords. In *Symposium on Usable Security and Privacy*, 2006.
- [55] M. D. Leonhard and V. Venkatakrisnan. A comparative study of three random password generators. In *IEEE International Conference on Electro/Information Technology*, 2007.
- [56] A. Léturgie. Un cas d’extragrammaticalité particulier: les amalgames lexicaux fantaisistes. *Linguistica*, 51: 87–104, 2011.
- [57] M. P. Lewis, G. F. Simons, and C. D. Fennig. *Ethnologue: Languages of the world*, 2014. Retrieved October 28, 2014 from <http://www.ethnologue.com/statistics/size>.
- [58] P. A. Luce. *Neighborhoods of words in the mental lexicon*. PhD thesis, Indiana University, 1986.
- [59] S. Majerus, M. Van der Linden, L. Mulder, T. Meulmans, and F. Peters. Verbal short-term memory reflects the sublexical organization of the phonological language network: evidence from an incidental phonotactic learning paradigm. *Journal of Memory and Language*, 51:297–306, 2004.
- [60] M. H. Messer, P. P. M. Leserman, J. Boom, and A. Y. Mayo. Phonotactic probability effect in nonword recall and its relationship with vocabulary in monolingual and bilingual preschoolers. *Journal of Experimental Child Psychology*, 105(4):306–323, 2010.
- [61] J. H. Neely. Semantic priming effects in visual word recognition: a selective review of current findings and theories. In *Basic processes in reading: visual word recognition*, pages 264–336. 1991.
- [62] NIST. Automated Password Generator (APG). Technical report, 1993.
- [63] M. Ohala. The syllable in Hindi. In H. van der Hulst and N. Ritter, editors, *The syllable: views and facts*,

- chapter 5, pages 93–111. Walter de Gruyter, 1999.
- [64] C. E. Piñeros. The creation of portmanteaus in the extragrammatical morphology of Spanish. *Probus*, 16(2):203–240, 2004.
- [65] J. O. Pliam. On the incomparability of entropy and marginal guesswork in brute-force attacks. In *Progress in Cryptology - INDOCRYPT*, volume 1977 of *Lecture Notes in Computer Science*, pages 67–79. 2000.
- [66] S. N. Porter. A password extension for improved human factors. *Computers & Security*, 1(1):54–56, 1982.
- [67] A. Rabkin. Personal knowledge questions for fallback authentication. In *Symposium on Usable Security and Privacy*, 2008.
- [68] A. Ralli and G. J. Xydopoulos. Blend formation in Modern Greek. In V. Renner, F. Maniez, and P. J. L. Arnaud, editors, *Cross-disciplinary perspectives on lexical blending*, pages 35–50. de Gruyter Mouton, Berlin, 2012.
- [69] E. Ronneberger-Sibold. Blending between grammar and universal cognitive principles: evidence from German, Farsi, and Chinese. In V. Renner, F. Maniez, and P. Arnaud, editors, *Cross-disciplinary perspectives on lexical blending*, number 252 in *Trends in Linguistics/Studies and monographs*, pages 115–143. de Gruyter, 2012.
- [70] S. Schechter, A. B. Brush, and S. Egelman. It’s no secret. measuring the security and reliability of authentication via ‘secret’ questions. In *IEEE Symposium on Security & Privacy*, 2009.
- [71] I. W. Schmidt, I. J. Berg, and B. G. Deelman. Relations between subjective evaluations of memory and objective memory performance. *Perceptual and Motor Skills*, 93(3), 2001.
- [72] K. E. Shaw. Head faithfulness in lexical blends: A positional approach to blend formation. Master’s thesis, 2013.
- [73] K. E. Shaw, A. M. White, E. Moreton, and F. Monrose. Emergent faithfulness to morphological and semantic heads in lexical blends. In J. Kingston, C. Moore-Cantwell, J. Pater, and R. Staubs, editors, *Proceedings of 2013 Meetings on Phonology*, Washington, DC, 2014. Linguistic Society of America.
- [74] R. Shay, S. Komanduri, P. G. Kelley, P. G. Leon, M. L. Mazurek, L. Bauer, N. Christin, and L. F. Cranor. Encountering stronger password requirements. In *Symposium on Usable Security and Privacy*, 2010.
- [75] R. Shay, P. G. Kelley, S. Komanduri, M. L. Mazurek, B. Ur, T. Vidas, L. Bauer, N. Christin, and L. F. Cranor. Correct horse battery staple: Exploring the usability of system-assigned passphrases. In *Symposium on Usable Security and Privacy*, 2012.
- [76] S. L. Smith. Authenticating users by word association. *Computers & Security*, 6(6):464–470, 1987.
- [77] A. Somayaji, D. Mould, and C. Brown. Towards narrative authentication: or, against boring authentication. In *Workshop on New Security Paradigms*, 2013.
- [78] S. L. Thompson-Schill, K. J. Kurtz, and J. D. E. Gabrieli. Effects of semantic and associative relatedness on automatic priming. *Journal of Memory and Language*, 38:440–458, 1998.
- [79] A. S. C. Thorn and C. R. Frankish. Long-term knowledge effects on serial recall of nonwords are not exclusively lexical. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 31(4):729–735, 2005.
- [80] J. Thorpe, P. C. van Oorschot, and A. Somayaji. Pass-thoughts: authenticating with our minds. In *Workshop on New Security Paradigms*, 2005.
- [81] J. Thorpe, A. Salehi-Abari, and R. Burden. Video-passwords: advertising while authenticating. In *Workshop on New Security Paradigms*, 2012.
- [82] R. Treiman. The structure of spoken syllables: evidence from novel word games. *Cognition*, 15:49–74, 1983.
- [83] R. Treiman. The division between onsets and rimes in english syllables. *Journal of Memory and Language*, 25:476–491, 1986.
- [84] R. Treiman, B. Kessler, S. Knewasser, R. Tincoff, and M. Bowman. English speakers’ sensitivity to phonotactic patterns. In *Papers in Laboratory Phonology V: Acquisition and the Lexicon*, pages 269–282. 2000.
- [85] TV Tropes Foundation. Television tropes and idioms: Portmanteau couple name, 2012. Retrieved November 11, 2012 from <http://tvtropes.org/pmwiki/pmwiki.php/Main/PortmanteauCoupleName>.
- [86] R. Weide. The Carnegie Mellon pronouncing dictionary, 1998. URL <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>.
- [87] M. Weir, S. Aggarwal, M. Collins, and H. Stern. Testing metrics for password creation policies by attacking large sets of revealed passwords. In *ACM Conference on Computer and Communications Security*, 2010.
- [88] A. M. White, K. Z. Snow, A. Matthews, and F. Monrose. Hookt on fon-iks: Phonotactic Reconstruction of Encrypted VoIP Conversations. In *IEEE Symposium on Security and Privacy*, 2011.
- [89] Y. Zhang, F. Monrose, and M. K. Reiter. The security of modern password expiration. In *ACM Conference on Computer and Communications Security*, 2010.
- [90] P. Zimmermann, A. Johnston, and J. Callas. ZRTP: Media Path Key Agreement for Unicast Secure RTP. RFC 6189, 2011. URL <http://www.ietf.org/rfc/rfc6189.txt>.
- [91] M. Zviran and W. J. Haga. Cognitive passwords: The key to easy access control. *Computers & Security*, 9(8):723–736, 1990.
- [92] M. Zviran and W. J. Haga. Password security: an empirical study. *Management Information Systems*, 15(4):161–185, 1999.