

Guess Again (and Again and Again): Measuring Password Strength by Simulating Password-Cracking Algorithms

Patrick Gage Kelley, Saranga Komanduri, Michelle L. Mazurek, Rich Shay, Tim Vidas Lujo Bauer, Nicolas Christin, Lorrie Faith Cranor, Julio Lopez

August 31, 2011

CMU-CyLab-11-008

CyLab
Carnegie Mellon University
Pittsburgh, PA 15213

Guess again (and again and again): Measuring password strength by simulating password-cracking algorithms

Patrick Gage Kelley, Saranga Komanduri, Michelle L. Mazurek, Rich Shay, Tim Vidas
Lujó Bauer, Nicolas Christin, Lorrie Faith Cranor, Julio López

Carnegie Mellon University Pittsburgh, PA

{pgage,sarangak,mmazurek,rshay,tvidas,lbauer,nicolasc,lorrie,julio.lopez}@cmu.edu

Abstract

Text-based passwords remain the dominant authentication method in computer systems, despite significant advancement in attackers' capabilities to perform password cracking. In response to this threat, password composition policies have grown increasingly complex. However, there is insufficient research defining metrics to characterize password strength and evaluating password-composition policies using these metrics. In this paper, we describe an analysis of 12,000 passwords collected under seven composition policies via an online study. We develop an efficient distributed method for calculating how effectively several heuristic password-guessing algorithms guess passwords. Leveraging this method, we investigate (a) the resistance of passwords created under different conditions to password guessing; (b) the performance of guessing algorithms under different training sets; (c) the relationship between passwords explicitly created under a given composition policy and other passwords that happen to meet the same requirements; and (d) the relationship between guessability, as measured with password-cracking algorithms, and entropy estimates. We believe our findings advance understanding of both password-composition policies and metrics for quantifying password security.

1 Introduction

Text-based passwords are the most commonly used authentication method in computer systems. As shown by previous research (e.g., [2]), passwords are often easy for attackers to compromise. A common threat model is an attacker who steals a list of hashed passwords, enabling him to attempt to crack them offline at his leisure. The many recent examples of data breaches involving large numbers of hashed passwords (Booz Allen Hamilton, HBGary, Gawker, Sony Playstation, etc.), coupled with the availability of botnets that offer large computational resources to attackers, make such threats very real [3, 5, 6, 10]. Once these passwords have been cracked, they can be used to gain access not only to the original site, but also to other accounts where users have reused their passwords. This is an important consideration because studies indicate that password reuse (exactly and with minor variations) is a common and growing practice as users acquire more online accounts [18, 38].

To mitigate the danger of such attacks, system administrators specify password-composition policies. These policies force newly created passwords to adhere to various requirements intended to make them harder to guess. Typical requirements are that passwords include a number or a symbol, that they exceed a certain minimum length, and that they are not words found in a dictionary.

Although it is generally believed that password-composition policies make passwords harder to guess, and hence more secure, research has struggled to quantify the level of resistance to guessing provided by different password-composition policies or the individual requirements they comprise. The two most commonly used methods for quantifying the effect of password-composition policies are estimating the entropy of passwords induced by password-composition policies using NIST guidelines [8], and empirically analyzing passwords created under different password-composition policies with password-guessing tools (e.g., [46]). The former, however, is not based on empirical data, and the latter is difficult to apply because of the dearth of password sets created under different password-composition policies.

In this paper, we take a substantial step forward in understanding the effects of password-composition policies on the guessability of passwords. First, we compile a dataset of 12,000 plaintext passwords collected from different

participants under seven different password-composition policies during a six-month online study. Second, we develop approaches for calculating how long it would take for various password-guessing tools to guess each of the passwords we collected. This allows us to evaluate the impact on security of each password-composition policy.

Contributions. Our paper makes the following specific contributions:

1. We implement a distributed technique (*guess-number calculator*) to determine if and when a given password-guessing algorithm, trained with a given data set, would guess a specific password. This allows us to evaluate the effectiveness of password-guessing attacks much more quickly than we could using existing cracking techniques.
2. We compare, more accurately than was previously possible, the guessability of passwords created under different password-composition policies. Because of the efficiency of our calculations (compared to guessing passwords directly), we can investigate the effectiveness of multiple password-guessing approaches with multiple tunings. Our findings show that a password-composition policy requiring long passwords with no other restrictions provides (relative to other policies we tested) excellent resistance to guessing.
3. We study the impact of tuning on the effectiveness of password-guessing algorithms. We also investigate the significance of test-set selection when evaluating the strength of different password-composition policies.
4. We investigate the effectiveness of entropy as a measure of password guessability. For each composition policy, we compare our guessability calculations to two independent entropy estimates: one based on the NIST guidelines mentioned above, and a second that we calculate empirically from the plaintext passwords in our dataset. We find that both measures of entropy have only very limited relationships to password strength as measured by guessability.

Mechanical Turk and controlled password collection. As with any user study, it is important to reflect on the origin of our dataset to understand the generalizability of our findings. We collected a dataset of 12,000 plaintext passwords using Amazon’s Mechanical Turk crowdsourcing service (MTurk). Many researchers have examined the use of MTurk workers (Turkers) as participants in human-subjects research [7, 15, 21, 22, 23, 31]. About half of all Turkers are American, with Indian participation increasing rapidly in the last 2-3 years to become about one third of Turkers [31]. American Turkers are about two-thirds women, while Indian Turkers are similarly weighted toward men [21]. Overall, the Turker population is younger and more educated than the general population, with 40% holding at least a bachelor’s degree; both of these trends are more pronounced among Indian Turkers [21, 31].

Buhrmester et al. find that the Turker population is significantly more diverse than samples used in typical lab-based studies that heavily favor college-student participants [7]. This study, and others, found that well-designed MTurk tasks provide high-quality user-study data [7, 15, 23, 41].

This analysis of MTurk has important implications in the context of studying passwords. We expect our findings will be more generalizable than those from lab studies with a more constrained participant base. Because we collected demographic information from our participants, our sample (and any biases it introduces) can be more accurately characterized than samples based on stolen password lists from various websites collected under uncertain circumstances.

A related consideration is that while our participants created real passwords that were needed several days later to complete the study and obtain a small bonus payment, these passwords did not protect high-value accounts. Password research has consistently been limited by the difficulty of studying passwords used for high-value accounts. Lab studies have asked participants to create passwords that protect simulated accounts, \$5, a chance to win an iPod in a raffle, or access to university course materials including homework and grades [9, 11, 26, 49]. Other studies have relied on the leaked password lists like the RockYou set [42, 46]. While this set contains millions of passwords, it also contains non-password artifacts that are difficult to filter out definitively, its provenance and completeness are unclear, and it is hard to say how much value users place on protecting an account from a social gaming service. Other commonly used leaked password lists come from sites including MySpace, *silentwhisper.net*, and a variety of Finnish websites, with user valuations that are similarly difficult to assess [13, 47].

Overall, although our dataset is not ideal, we contend that our findings do provide significant insight into the effects of password-composition policies on password guessability. Because so little is known about this important topic, even imperfect information constitutes progress.

Roadmap. We proceed as follows. In Section 2 we survey related work. We describe our data collection and analysis methodology in Sections 3 and 4. We convey our main results in Section 5, and discuss the generalizability of our findings and some ethical considerations in Section 6. We conclude in Section 7 with a discussion of the applicability of our work to future research and the implications of our findings on defining practical password-composition policies.

2 Background and related work

Research on passwords has been active for many years. In this section, we review the work most closely related to our research, first summarizing the different types of data collection and analysis that have been used. We then discuss work focused on evaluating the impact of password policies, followed by metrics proposed to evaluate password strength.

Collection and analysis of password data. Many prior studies of passwords have used small sample sizes [19, 24, 29, 51], obtained through user surveys or lab studies. Kuo et al. asked 290 users to create passwords through an online survey and used a password-cracking tool to estimate the security of the passwords so acquired [26]. We also use an online survey, but we consider larger and more varied sets of passwords. In addition, we recruit participants using Mechanical Turk, which produces more diverse samples than typical lab studies [7].

Other studies analyze large samples of passwords ostensibly created by users for actual accounts of varying importance [2, 4, 13, 16, 46, 50]. Unlike these studies, we study the impact of different password policies on password strength and use passwords collected under controlled password-policy conditions.

Impact of password policies. Several studies have considered the impact of different password policies on password strength. In lab studies, Proctor et al. [30] and Vu et al. [44] found that passwords created under stricter composition requirements were more resistant to automated cracking, but also more difficult for participants to create and remember. We consider similar data, but for a much larger set of users, allowing us to evaluate the effectiveness of various requirements more comprehensively. Other findings suggest that too-strict policies (those that make creating and remembering passwords too difficult) induce coping strategies that can hurt both security and productivity [1, 20, 36, 37, 40]. Further, Florêncio and Herley found that the strictest policies are often used not by organizations with high-value assets to protect, but organizations that do not have to compete on customer service [17].

An increasingly popular password-strengthening measure that we also investigate is subjecting new passwords to a blacklist check. Spafford demonstrated that a blacklist check can be performed in constant time regardless of size [39], making large blacklist checking feasible. Schechter et al. proposed a password policy in which passwords chosen by too many users are blacklisted for subsequent users [32]. This offers many theoretical advantages over other password-composition schemes.

Measuring password strength. Effective evaluation of password strength requires defining a proper metric. One possible metric is information entropy, defined by Shannon as the expected value (in bits) of the information contained in a string [34]. Massey connects entropy with password strength by demonstrating that entropy provides a lower bound on the expected number of guesses to find a text [28]. A 2006 National Institute of Standards and Technology (NIST) publication uses entropy to represent the strength of a password [8]. Verheul derives a theoretical distribution of variable-length passwords with optimal entropy and guess resistance [43]. Neither calculated entropy empirically. Florêncio and Herley estimated theoretical entropy for the field data they analyzed [16].

An alternative to entropy as a metric of password strength is the notion of “guessability,” which characterizes the time needed by an efficient password-cracking algorithm to discover a password. In one use of this metric, Weir et al. divide a large set of existing passwords into different categories based on composition, then apply automated cracking tools to examine how well NIST’s entropy estimates predict measured guessing difficulty [46]. Similarly to our work, Dell’Amico et al. [13] also attempt to evaluate password strength by calculating guessing probabilities yielded by popular password-cracking heuristics.

Narayanan et al. discuss a password-cracking technique based on a Markov model, in which password guesses are made based on contextual frequency of characters [29]. Marechal [27] and Weir [45] both examine this model and find it more effective for password cracking than the popular password-cracking program John the Ripper [14]. Weir et al. present a novel password-cracking technique that uses the text structure from training data while applying mangling rules to the text itself [47]. The authors found their technique to be more effective than John the Ripper. In a separate study, Zhang et al. found Weir’s algorithm most effective among the techniques they used [50].

In this work, we apply the Weir et al. algorithm and a variation of the Markov model to generate blacklists that restrict password creation in some of our study conditions, and as the basis for one implementation of a new measure of password strength, the *guess number*, which we apply to user-created passwords collected under controlled password-composition policies.

3 Methodology: Data collection

In this section, we discuss our methodology for collecting plaintext passwords, the word lists we used to assemble the blacklists used in some conditions, and the eight conditions under which we gathered data. We also summarize participant demographics.

3.1 Collection instrument

From August 2010 to January 2011, we advertised a two-part study on Mechanical Turk, paying between 25 and 55 cents for the first part and between 50 and 70 cents for the second part. The consent form indicated the study pertained to visiting secure websites.

Each participant was given a scenario for making a new password, then asked to create a password that met a set of password-composition requirements; the scenarios and requirements are detailed in Section 3.3. Participants who entered a password that did not conform to requirements were shown an error message indicating which requirements were not met, then asked to try again until a satisfactory password was created. After creating a password, participants took a brief survey about demographics and password creation. Participants were then asked to recall the password just created; after five failed attempts, the password was displayed. For the second part of the study, participants were emailed two days later and asked to return to the website and recall their passwords. We also measured the incidence of passwords being written down or otherwise stored (via detecting browser storage and copy-paste behavior, as well as asking participants; see Section 6 for details). Unless otherwise noted, only data from the first part of the study is reported in this paper. Data from the second part of the study, which is primarily used to assess memorability and usability factors, is omitted due to space constraints. Prior research has considered memorability and usability factors for a subset of the policies we examine [25]; we briefly revisit these findings when we discuss our results in Section 5.

3.2 Word lists for algorithm training

We use six publicly available word lists as training data in our analysis and to assemble the blacklists used in some of our experimental conditions. The *RockYou* password set [42] includes more than 30 million passwords, and the *MySpace* password set [33] contains about 45,000 passwords. (We discuss ethical considerations related to these datasets in Section 6.) The *inflection list*¹ contains words in varied grammatical forms such as plurals and past tense. The *simple dictionary* contains about 200,000 words and is a standard English dictionary available on most Unix systems. We also used two cracking dictionaries from the Openwall Project² containing standard and mangled versions of dictionary words and common passwords. The *free Openwall list* contains about 4 million words, while the *paid Openwall list* contains more than 40 million. While these data sources are not ideal, they are publicly available; we expect attackers would use these word lists or others like them for training data. In Section 5.2, we consider the effect of a variety of training sets drawn from these word lists as well as our collected password data.

3.3 Conditions

Our participants were divided into eight conditions comprising seven sets of password-composition requirements and two password-creation scenarios. We used two scenarios in order to measure the extent to which giving participants different instructions affects password strength. The *survey scenario* was designed to simulate a scenario in which users create low-value passwords, while the *email scenario* was designed to elicit higher-value passwords. All but one condition used the email scenario.

In the *survey scenario*, participants were told, “To link your survey responses, we will use a password that you create below; therefore it is important that you remember your password.”

In the *email scenario*, participants were told, “Imagine that your main email service provider has been attacked, and your account became compromised. You need to create a new password for your email account, since your old password may be known by the attackers. Because of the attack, your email service provider is also changing its password rules. Please follow the instructions below to create a new password for your email account. We will ask you to use this password in a few days to log in again, so it is important that you remember your new password. Please take the steps you would normally take to remember your email password and protect this password as you normally would protect the password for your email account. Please behave as you would if this were your real password!”

¹<http://wordlist.sourceforge.net>

²<http://www.openwall.com/wordlists/>

The eight conditions are detailed below.

basic8survey: Participants were given the survey scenario and the password-composition policy “Password must have at least 8 characters.” This is the only condition using the survey scenario.

basic8: Participants were given the email scenario and the password-composition policy “Password must have at least 8 characters.” Only the scenario differentiates this from basic8survey.

basic16: Participants were given the email scenario and the password-composition policy “Password must have at least 16 characters.”

dictionary8: Participants were given the email scenario and the password-composition policy “Password must have at least 8 characters. It may not contain a dictionary word.” We performed a dictionary check by removing non-alphabetic characters and checking the remainder against a dictionary, ignoring case. This method is used in practice, including at our institution. We used the free Openwall list as the dictionary.

comprehensive8: Participants were given the email scenario and the password-composition policy “Password must have at least 8 characters including an uppercase and lowercase letter, a symbol, and a digit. It may not contain a dictionary word.” We performed the same dictionary check as in dictionary8. This condition reproduced NIST’s comprehensive password-composition requirements [8].

blacklistEasy: Participants were given the email scenario and the password-composition policy “Password must have at least 8 characters. It may not contain a dictionary word.” We checked the password against the simple Unix dictionary, ignoring case. Unlike the dictionary8 and comprehensive8 conditions, the password was not stripped of non-alphabetic characters before the check.

blacklistMedium: This condition is the same as the blacklistEasy condition, except we used the paid Openwall list.

blacklistHard: This condition is the same as the blacklistEasy condition, except we used a five-billion-word dictionary we created using the algorithm outlined by Weir et al. [47]. For this condition, we trained Weir et al.’s algorithm on the MySpace, RockYou, and inflection lists. Both the training and testing were conducted case-insensitively, increasing the strength of the blacklist.

3.4 Participant demographics

Among participants who completed part one of our study, 55% returned within 3 days and completed part two. We detected no statistically significant differences in the guessability of passwords between participants who took just the first part of the study and those who participated in both parts. As a result, to maximize the participant data in our analyses and use the same number of participants for each condition, our dataset includes passwords from the first 1,000 participants in each condition to successfully complete the first part of the study. To conduct a wider variety of experiments, we used data from an additional 2,000 participants each in basic8 and comprehensive8.

Among these 12,000 participants, 53% percent reported being male and 45% female, with a mean reported age of 29 years. This makes our sample more male and slightly younger than Mechanical Turk participants in general [7,31]. About one third of participants reported studying or working in computer science or a related field. The proportion related to computer science did not vary significantly across conditions, except between blacklistEasy and blacklistHard (38% to 31%, respectively; pairwise Holm-corrected Fisher’s exact test [PHFET], $p < 0.03$). Participants in the basic16 condition were slightly but significantly older (mean 30.3 years) than those in blacklistHard, basic8, and comprehensive8 (means 28.6, 28.9, and 29.1 years respectively; PHFET, $p < 0.03$). We observed no significant difference in gender between any pair of conditions (PHFET, $p > 0.05$).

4 Methodology: Data analysis

This section explains how we analyzed our collected password data. First, and most importantly, Section 4.1 discusses our approach to measuring how resistant passwords are to cracking, i.e., guessing by an adversary. We present a novel, efficient method that allows a broader exploration of guessability than would otherwise be possible. For comparison purposes, we also compute two independent entropy approximations for each condition in our dataset, using methods described in Section 4.2.

4.1 Guess-number calculators

Traditionally, password guess resistance is measured by running one or more password-cracking tools against a password set and recording when each password is cracked. This works well when the exploration is limited to a relatively

small number of guesses (e.g., 10^{10} , or roughly the number of guesses a modern computer could try in one day). However, as the computational power of potential adversaries increases, it becomes important to understand how many passwords an adversary could crack with many more guesses.

To this end, we introduce the *guess number calculator*, a novel method for measuring guess resistance more efficiently. We take advantage of the fact that, for most deterministic password-guessing algorithms, it is possible to create a calculator function that maps a password to the number of guesses required to guess that password. We call this output value the *guess number* of the password. A new guess number calculator must be implemented for each cracking algorithm under consideration. For algorithms like [46] that use a *training set* of known passwords to establish guessing priority, a new *tuning* of the calculator is generated for each new training set to be tested.

Because we collect plaintext passwords, we can use a guessing algorithm’s calculator function to look up the associated guess number for each password, without actually running the algorithm. This works for the common case of deterministic guessing algorithms (e.g., [14, 27, 29, 46]).

We use this approach to measure the guessability of a set of passwords in several ways. We compute the percentage of passwords that would be cracked by a given algorithm, which is important because the most efficient cracking tools use heuristics and do not explore all possible passwords. We can also compute the percentage that would be cracked within a given number of guesses, or the number of guesses required to crack a certain percentage of passwords. We also use calculators to compare the performance of different cracking algorithms, and different training-set tunings within each algorithm. By combining guess-number results across a variety of algorithms and training sets, we can develop a general picture of the overall strength of a set of passwords.

We implemented two guess-number calculators: one for a brute-force algorithm loosely based on the Markov model, and one for the heuristic algorithm proposed by Weir et al., which is currently the state-of-the-art approach to password cracking [46, 50]. We selected these two algorithms as the most promising brute-force and heuristic options, respectively, after comparing the passwords we collected to lists of 1, 5, and 10 billion guesses produced by running a variety of cracking tools and tunings. From this point forward, we will refer to them as the brute-force Markov (BFM) and Weir algorithms.

4.1.1 Training sets

Both algorithms for which we implemented calculators require a *training set*: a corpus of known passwords used to generate a list of guesses and determine in what order they should be tried.

We explore a varied space of training sets constructed from different combinations of the publicly available word lists described in Section 3.2 and subsets of the passwords we collected. This allows us to assess whether complementing publicly available data with passwords collected from the system under attack improves the performance of the cracking algorithms. We further consider training-set variations specifically tuned to our two most complex policy conditions, comprehensive8 and basic16.

Each of our experiments calculates guess numbers only for those passwords on which we did not train, using a cross-validation approach. For a given experiment, we split our passwords into n partitions, or *folds*. We generate a training set from public data plus $(n - 1)$ folds of our data, and test it on the remaining fold. We use each of the n folds as test data exactly once, requiring n iterations of testing and training. We recombine results from the n folds, yielding guess-number results for all of our passwords. Because training often involves significant computational resources, as described in Section 4.1.3, we limit to two or three the number of iterations in our validation. Based on the similarity of results we observed between iterations, this seems sufficient. We describe our training and test sets in detail in Appendix A.

We do not claim these training sets or algorithms represent the optimal technique for guessing the passwords we collected; rather, we focus on comparing guess resistance across password-composition policies. Investigating the performance of guessing algorithms with different tunings also provides insight into the kind of data set an attacker might need in order to efficiently guess passwords created under a specific password-composition policy.

4.1.2 BFM calculator

The BFM calculator determines guess numbers for a brute-force cracking algorithm loosely based on Markov chains [27, 29]. Our algorithm differs from previous work by starting with the minimum length of the password policy, and increasing the length of guesses until all passwords are guessed. Unlike other implementations, this covers the entire password space, but does not try guesses in strict probability order.

The BFM algorithm uses the training set to calculate the frequency of first characters and of digrams within the password body, and uses these frequency to deterministically construct guessing orders of unknown passwords. For

example, assume an alphabet of $\{A, B, C\}$ and a three-character-minimum configuration. If A is the most likely starting character learned from the training data, the character most likely to follow A is B , and the character most likely to follow B is C , then the first guess will be ABC . If the next-most-likely character to follow B is A , the second guess will be ABA , and so forth.

Our guess-number calculator for this algorithm processes the training data to generate a lookup table that maps each string to the number of guesses needed to reach it, as follows. For an alphabet of N characters, and passwords of length L , any time the first character tried is incorrect, we know that the algorithm will try N^{L-1} incorrect guesses before switching to a different first character. So, if the first character of the password to be guessed is the k -th character to be tried, there will be at least $(k-1)N^{L-1}$ incorrect guesses. We can then iterate the computation: when the first character is correct, but the second character is incorrect, the algorithm will try N^{L-2} incorrect guesses, and so forth. By looking up the order in which characters are tried, we can then simply add up the total number of incorrect guesses to discover how many iterations will be needed before hitting a successful guess for a given password, without having to actually try the guesses.

4.1.3 Weir algorithm calculator

We also apply the principle of calculating guess numbers to Weir et al.’s more complex algorithm, both to demonstrate the feasibility of our approach and to evaluate the strength of our collected password sets. The Weir algorithm is explained in detail in [47], and uses the following definitions: *structures* are patterns of character types such as letters, digits, and symbols; a *terminal* is one instantiation of a structure; and a *probability group* is a set of terminals with the same probability of occurring.

As with the BFM calculator, we process training data to create a lookup table, then calculate the guess number for each password. The mechanism for processing training data is outlined in Algorithm 1. To calculate the guess number for a password, we determine that password’s probability group. Using the lookup table created from the training set, we determine how many guesses would be required to reach that probability group. We then add the number of guesses required to reach the exact password within that probability group. This works because once the Weir algorithm reaches a given probability group, all terminals in that group are tried in a deterministic order.

Because creating this lookup table is time-intensive, we set a cutoff point—50 trillion guesses, which allows most Weir-calculator experiments to run in 24 hours or less in our setup—past which we do not calculate the guess number for additional passwords. By checking the available terminals, we can still determine whether passwords that are not guessed by this point will ever be guessed, but not exactly when they will be guessed.

Algorithm 1 Creation of a lookup table which, given a probability group, returns the number of guesses required for the Weir algorithm to begin guessing terminals of that group. An *l.c.s.* is a *longest common substring*, the longest substrings in a probability group made from characters of the same type. For example, for $UUss9UUU$, the *l.c.s.*’s would be UU , ss , 9 , and UUU . (In this example, U represents uppercase letters, s represents lowercase letters, and 9 represents digits.)

```

 $\mathcal{T}$  = New Lookup Table
for all structures  $s$  do
  for all probability_group  $pg \in s$  do
    for all l.c.s.  $\in pg$  do
       $c_i$  = Number of terminals of l.c.s.
       $p_i$  = Probability of l.c.s. in training data
    end for
     $probability = \prod p_i$ 
     $\mathcal{T}.add: probability, pg, \prod c_i$ 
  end for
end for
Sort( $\mathcal{T}$ ) by probability
Add to each value in ( $\mathcal{T}$ ) the sum of prior values

```

Distributed computation. Calculating guess numbers for Weir’s algorithm becomes data intensive for the sets of structures used in this work. More specifically, Algorithm 1 generates a large number of elements to build the lookup table \mathcal{T} . To accelerate the process, we implemented a distributed version of Algorithm 1 as follows. We split the top-most loop into coarse-grained units of work that are assigned to m tasks, each of which processes a subset of the

structures in s . Each task reads a shared dictionary with the training data and executes the two internal loops of the algorithm. Each iteration of the loop for the probability groups in s emits an intermediate tuple. The intermediate tuples produced by all the tasks are grouped by probability ranges, sorted, and stored. A final sequential pass over the sorted table adds the sum of prior values.

We implemented our distributed approach using Hadoop [48], an open-source version of the MapReduce framework [12]. In the implemented approach, while all m tasks receive equally sized subsets of the input, the tasks perform different amounts of work depending on the complexity of the structures in each respective input subset. As a result, task execution times vary widely. Nevertheless, this approach enabled us to compute guess numbers for sets of 4000 passwords in hours, rather than days, in a 64-node Hadoop cluster. The resulting lookup tables store on the order of hundreds of billions of elements with their associated probabilities and occupy up to 1.3 TB of storage each.

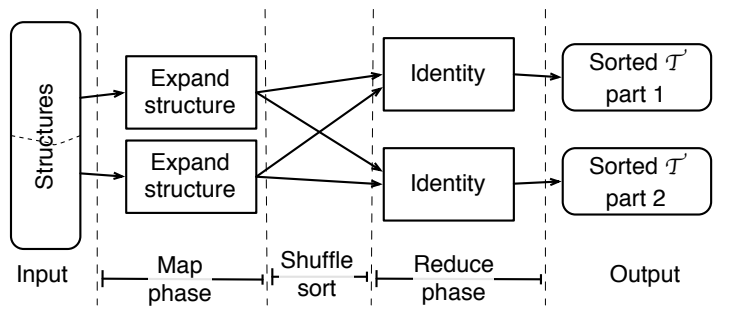


Figure 1: Distributed computation of the Weir-algorithm guess numbers using Hadoop. The framework splits the set of input structures and executes a map task for each subset. The map tasks emit tuples of the form $\langle probability, pg, \Pi c_i \rangle$. The framework sorts the the intermediate tuples using a *total ordering partitioner*. A reduce task, which in this case is the identity reducer, writes to storage a subset of the sorted tuples.

4.2 Entropy

In order to investigate how well entropy estimates correlate with guess resistance in practice, we compare our guess number results for each condition in our dataset to two independently calculated entropy approximations. First, we apply the commonly used NIST guidelines, which suggest that each password-composition rule contributes a specific amount of entropy and that the entropy of the policy is the sum of the entropy contributed by each rule. Our second approximation is calculated empirically from the plaintext passwords in our dataset, using the technique described by Shay et al. [38]. In this method, we calculate for each password in the condition the entropy contributed by the number, content, and type of each character, using Shannon’s formula [35]. We then sum the individual entropy contributions to estimate the total entropy of the passwords in that condition.

5 Findings

We calculated guess numbers under 32 different combinations of algorithm and training data. Although we do not have space to include all the raw results, we distill from them four major findings with application both to selecting password policies and to conducting password research:

- Among conditions we tested, basic16 provides the greatest security against a powerful attacker, outperforming the more complicated comprehensive8 condition. We also detail a number of other findings about the relative difficulty of cracking for the different password-composition policies we tested.
- Access to abundant, closely matched training data is important for successfully cracking passwords from stronger composition policies. While adding more and better training data provides little to no benefit against passwords from weaker conditions, it provides a significant boost against stronger ones.
- Passwords created under a specific composition policy do not have the same guess resistance as passwords selected from a different group that happen to meet the rules of that policy; effectively evaluating the strength of a password policy requires examining data collected under that policy.
- While a limited relationship between Shannon information entropy (computed and estimated as described in Section 4.2) and guessability can be observed, especially when considering attacks on the order of a trillion guesses or more, entropy can provide no more than a very rough approximation of overall password strength.

We discuss these findings in detail in the rest of this section. We introduce individual experiments in detail before discussing their results. For convenience, after introducing an experiment we may refer to it using a shorthand name

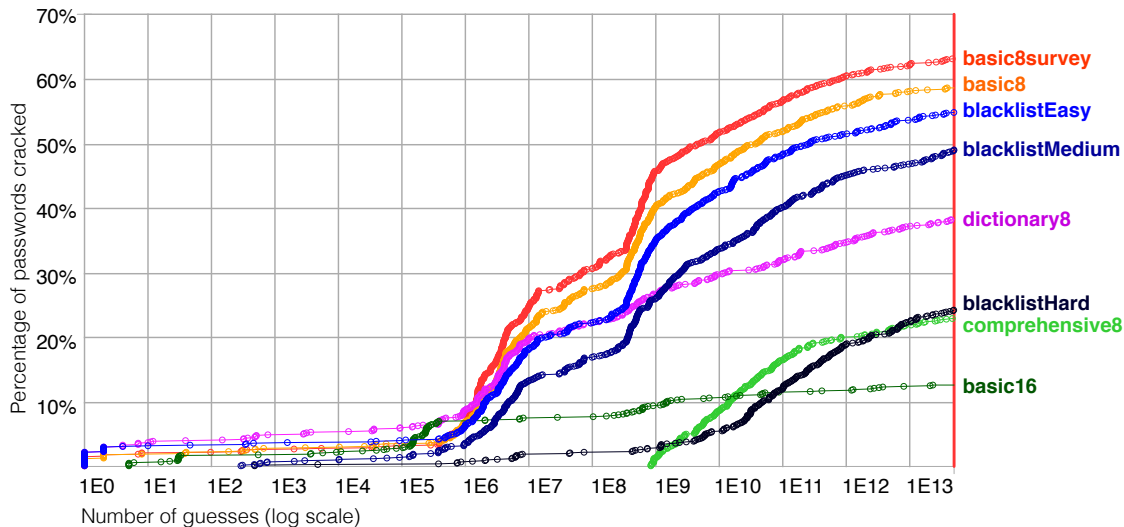


Figure 2: The number of passwords cracked vs. number of guesses, per condition, for experiment E. This experiment uses the Weir calculator and our most comprehensive training set, which combines our passwords with public data.

that maps to some information about that experiment, such as P for trained with public data, E for trained with everything, S for tested on password subsets, C8 for specialized training for comprehensive8, etc. A complete list of experiments and abbreviations can be found in Appendix A.

5.1 Comparing policies for guessability

In this section, we compare the guessability of passwords created under the seven password-composition policies we tested. We focus on two experiments that we consider most comprehensive. In each experiment we evaluate the guessability of all seven password-composition policies, but against differently trained guessing algorithms.

Experiment P4 is designed to simulate an attacker with access to a broad variety of publicly available data for training. It consists of a Weir-algorithm calculator trained on all the public word lists we use and tested on 1000 passwords from each condition. Experiment E simulates a powerful attacker with extraordinary insight into the password sets under consideration. It consists of a Weir-algorithm calculator trained with all the public data used in P4 plus 500 passwords from each of our eight conditions. We test on 500 other passwords from those conditions, with two-fold cross-validation for a total of 1000 test passwords. The results from experiments E and P4 are shown in Figures 2 and 3, respectively.

As suggested by these figures, which password-composition policy is best at resisting guessing attacks depends on how many attempts an attacker is expected to make. At one million and one billion guesses in both experiments, significantly fewer blacklistHard and comprehensive8 passwords were guessed than in any other condition.³ At one billion guesses in experiment E, 9.5%, 1.4%, and 2.9% of passwords were cracked in basic16, comprehensive8, and blacklistHard respectively; 40.3% of basic8 passwords were cracked.

As the number of guesses increases, basic16 begins to outperform the other conditions. At one trillion guesses, significantly fewer basic16 passwords were cracked than comprehensive8 passwords, which were in turn cracked significantly less than any other condition. After exhausting the Weir-algorithm guessing space in both experiments, basic16 remains significantly hardest to crack. The next best at resisting cracking were comprehensive8 and blacklistHard, performing significantly better than any of the other conditions. Condition comprehensive8 was significantly better than blacklistHard in experiment P4 but not in experiment E. In experiment E, 14.6, 26.4, and 31.0% of passwords were cracked in basic16, comprehensive8, and blacklistHard respectively; in contrast, 63.0% basic8 passwords were cracked.

Although guessing with the Weir algorithm proved more effective, we also compared the conditions using BFM. The findings (shown in Figure 4) are generally consistent with those discussed above: basic16 performs better than the other conditions.

³All comparisons in Sections 5.1, 5.2, and 5.3 tested using PHFET, significance level $\alpha = 0.05$.

Prior research examining the memorability and usability of a subset of the composition policies we examine here found that while in general less secure policies are more usable, basic16 is more usable than comprehensive8 by many measures [25]. This suggests basic16 is an overall better choice than comprehensive8.

It is important to note that 16-character-minimum policies are rare in practice. Hence, current guessing algorithms, including the Weir algorithm, are not built specifically with them in mind. Although we do not believe this affects our overall findings, further investigation would be beneficial.

5.2 Effects of training-data selection

Most practical cracking algorithms, including the ones we use, rely on training data to produce an ordering of guesses. As a result, it is important to consider how the choice of training data affects the success of password guessing, and consequently the guess resistance of a set of passwords. To address this, we examine the effect of varying the amount and source of training data on both total cracking success and on cracking efficiency. Interestingly, we find that the choice of training data affects different password-policy conditions differently; abundant, closely matched training data is critical when cracking passwords from harder-to-guess conditions, but less so when cracking passwords from easier ones.

For purposes of examining the impact of training data, the password-policy conditions we consider divide fairly neatly into two groups. For the rest of this section, we will refer to the harder-to-guess conditions of comprehensive8, basic16, and blacklistHard as *group 1*, and the rest as *group 2*.

Training with general-purpose data.

We first measure, via three experiments, the effect of increasing the amount and variety of training data. Experiment P3 was trained on public data including the MySpace and RockYou password lists as well as the inflection list and simple dictionary, and tested on 1000 passwords from each of our eight conditions. Experiment P4, as detailed in Section 5.1, was trained on everything from P3 plus the paid Openwall list. Experiment E, also described in 5.1, was trained on all the public data from P4 as well as 500 passwords from each of our conditions, using two-fold cross-validation. Figure 5 shows how these three training sets affect four example conditions, two from each group.

The cracking totals in each experiment reflect the overall increase in knowledge as training data is added. For group 1, adding Openwall increases total cracking on average 45%, while adding both Openwall and our data provides an average 96% improvement (all significant). In group 2, by contrast, the increases are more modest and only occasionally significant.

At one trillion and one billion guesses, the results are less straightforward, but increasing training data remains generally more effective against group 1 than group 2. Adding Openwall alone is not particularly helpful for group 1, providing few significant improvements at either guessing point, but it actually decreases cracking at one billion guesses significantly for several group 2 conditions. (We hypothesize this decrease occurs because Openwall is a

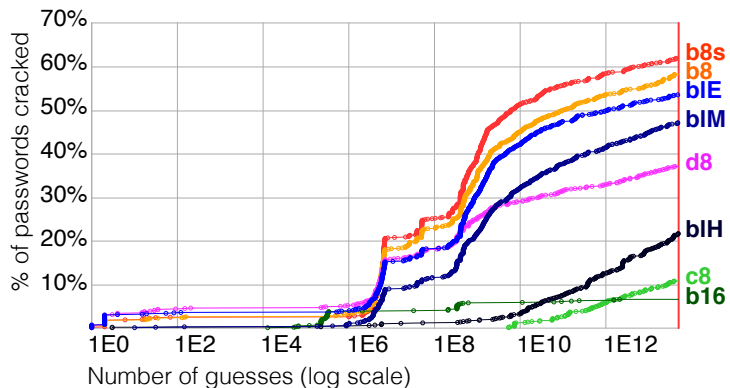


Figure 3: The number of passwords cracked vs. the number of guesses, per condition, for experiment P4. This experiment uses the Weir calculator and trains on a variety of publicly available data.

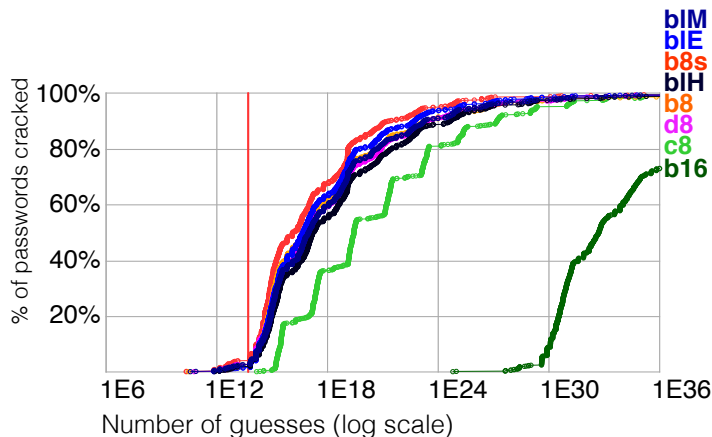


Figure 4: The number of passwords cracked vs. the number of guesses, using the BFM calculator trained on a combination of our data and public data (B2). We place a red vertical line at 50 trillion guesses to facilitate comparison with the Weir experiments. We stopped the Weir calculator at this point (as described in Section 4.1.3), but because the BFM algorithm is so much less efficient, we ran it for many more guesses in order to collect useful data.

dictionary and not a password set, so it adds knowledge of structures and strings at the cost of accurately assessing their probabilities.) At these guessing points, adding our data is considerably more effective for group 1 than adding Openwall alone, increasing cracking for each of the three conditions by at least 50% (all significant). By contrast, adding our data provides little to no improvement against group 2 conditions at either guessing point.

Taken together, these results demonstrate that increasing the amount and variety of information available in the training data provides significant improvement in cracking the harder-to-guess conditions, while providing little benefit and sometimes decreasing efficiency for the easier-to-guess conditions.

Training with specialized data. Having determined that training with specialized data is extremely valuable for cracking group 1 passwords, we wanted to examine what quantity of closely related training data is needed to effectively crack these “hard” conditions. For these tests, we focus on comprehensive8 as an example of a harder-to-guess condition, using the easier-to-guess basic8 condition as a control; for each of these conditions, we collected 3000 passwords.

We conducted five Weir-algorithm experiments, C8a through C8e, in which we trained on all the word lists described in Section 3.2, as well as between 500 and 2500 comprehensive8 passwords, in 500-password increments. For each experiment, we tested on the remaining comprehensive8 passwords. We also carried out a similar set of five experiments, B8a through B8e, in which we trained and tested with basic8 rather than comprehensive8 passwords.

Our results, illustrated in Figure 6, show that incrementally adding more of our collected data to the training set improves total cracking slightly for comprehensive8 passwords, but not for basic8. On average, for each 500 comprehensive8 passwords added to the training set, 2% fewer passwords remain uncracked. This effect is not linear, however; the benefit of additional training data levels off sharply between 2000 and 2500 training passwords. The differences between experiments begin to show significance around one trillion guesses, and increase as we approach the total number cracked.

For basic8, by contrast, adding more collected passwords to the training set has no significant effect on total cracking, with between 61 and 62% of passwords cracked in each experiment. No significant effect is observed at earlier guessing points including one million, one billion, or one trillion guesses, either.

One way to interpret this result is to consider the diversity of structures found in our basic8 and comprehensive8 password sets. The comprehensive8 passwords are considerably more diverse, with 1598 structures among 3000 passwords, as compared to only 733 structures for basic8. For comprehensive8, the single most common structure maps to 67 passwords, the most common 180 structures account for half of all passwords, and 1337 passwords have structures that are unique within the password set. By contrast, the most common structure in basic8 maps to 293 passwords, the top 13 structures account for half the passwords, and only 565 passwords have unique structures. As a result, small amounts of training data go considerably farther in cracking basic8 passwords than in comprehensive8.

Weighting training data. The publicly available word lists we used for training are all considerably larger than the number of passwords we collected. As a result, we needed to weight our data (i.e., include multiple copies in the training set) if we wanted it to have significant impact on the probabilities used by our guess-number calculators. Different weightings have no effect on the total number of passwords cracked, as all the same passwords are eventually guessed; however, they can affect the order and, therefore, the efficiency of guessing.

We tested three weightings, using 500 passwords from each of our eight conditions weighted to one-tenth, equal, and ten times the cumulative size of the included public lists. In each case, we tested on 500 other passwords from

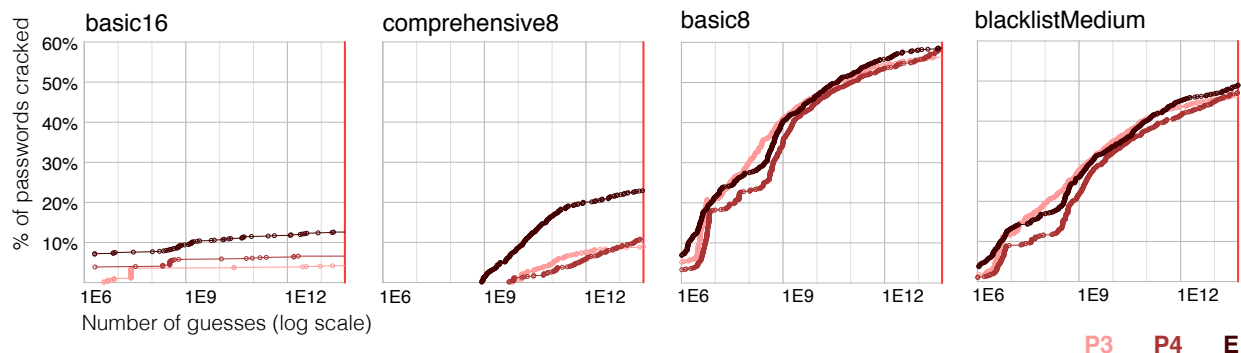


Figure 5: The effect of increasing training data by adding the Openwall list (experiment P4) and then our collected passwords (experiment E) on the effectiveness of cracking passwords, for four example conditions. Adding training data proves more helpful for the group 1 conditions (left two graphs) than for the others (right two graphs).

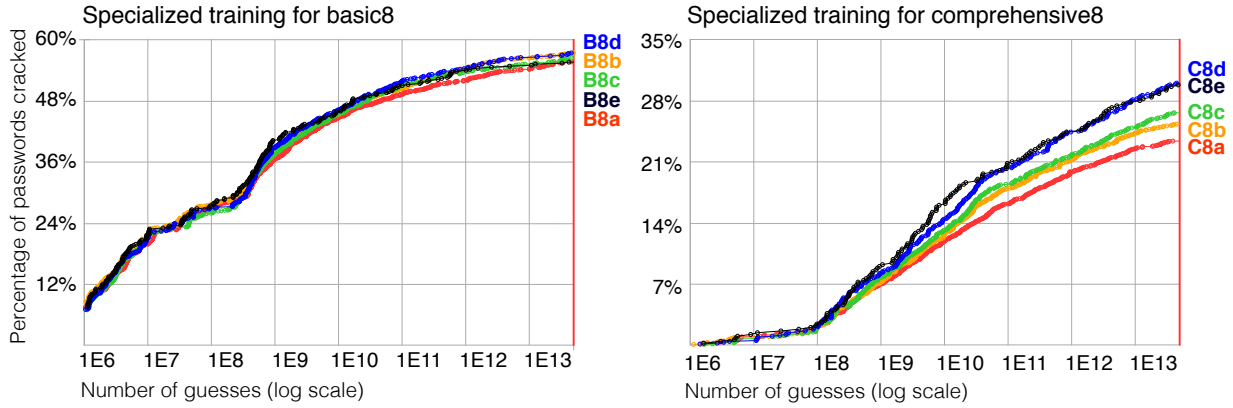


Figure 6: Left: Results of experiments B8a through B8e. Increasing the amount of specialized training data has limited effect on the basic8 condition. Right: Results of experiments C8a through C8e. Increasing the amount of specialized training data has a small but significant effect on the comprehensive8 condition.

each condition.

Overall, we found that weighting had only a minor effect. There were few significant differences at one million, one billion, or one trillion guesses, with equal weighting occasionally outperforming the other two in some conditions. From these results, we concluded that the choice of weighting was not particularly important, but we use an equal weighting in all other experiments that train with passwords from our dataset because it provides an occasional benefit.

BFM training. We also investigated the effect of training data on the performance of the BFM calculator, using four training sets: one with public data only, one that combined public data with collected passwords across our conditions, and one each specialized for basic8 and comprehensive8. Because the BFM algorithm eventually guesses every password, we were concerned only with efficiency, not total cracking. We found that adding our data had essentially no effect at either smaller or larger numbers of guesses. Specialized training for basic8 was similarly unhelpful. Specialized training for comprehensive8 does increase efficiency somewhat, reaching 50% cracked with about 30% fewer guesses.

5.3 Effects of test-data selection

Researchers typically don't have access to passwords created under the password-composition policy they want to study. To compensate, they start with a larger set of passwords (e.g., the RockYou set), and pare it down by discarding passwords that don't meet the desired composition policy (e.g., [16, 46]). A critical question, then, is whether subsets like these are representative of passwords actually created under a specific policy. We find that such subsets are not representative, and may in fact contain passwords that are more difficult to guess than passwords created under the policy in question.

In our experiments, we compared the guessability of 1000 comprehensive8 passwords to the guessability of the 206 passwords that meet the comprehensive8 requirements but were collected across our other seven conditions (the

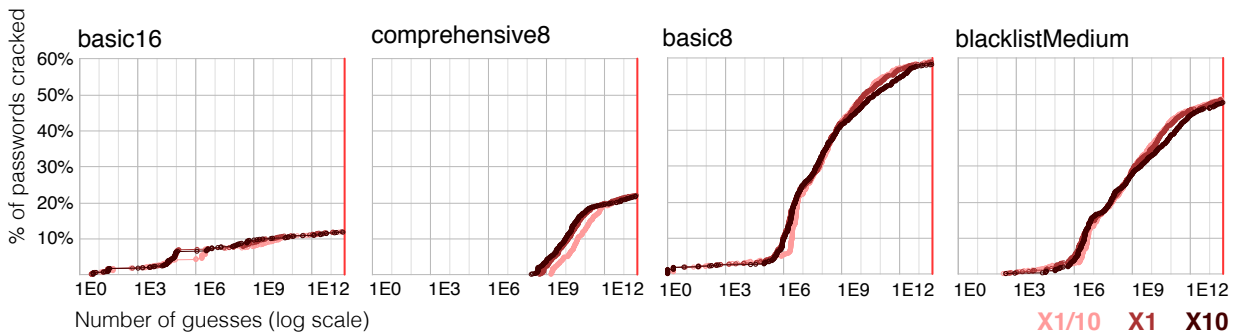


Figure 7: Varying the weighting of our passwords within the public training data among one-tenth (X1/10), equal weighting (X1), and ten times (X10) has little to no effect on the efficiency of cracking passwords. Results shown for four example conditions.

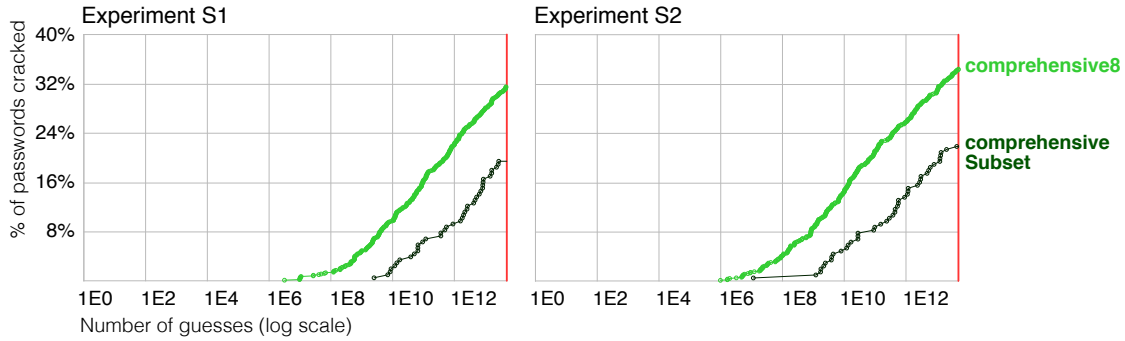


Figure 8: Passwords generated under the comprehensive8 condition proved significantly easier to guess than passwords that conform to the comprehensive8 requirements but are generated under other composition policies. In experiment S1, shown at left, the Weir calculator was trained with only public data; in experiment S2, shown at right, the Weir calculator was trained on a combination of our data and public data.

comprehensiveSubset set). We performed this comparison with two different training sets: public data, with an emphasis on RockYou passwords that meet comprehensive8 requirements (experiment S1); and the same data enhanced with our other 2000 collected comprehensive8 passwords (experiment S2).

Both experiments show significant differences between the guessability of comprehensive8 and comprehensiveSubset test sets, as shown in Figure 8. In the two experiments, 40.9% of comprehensive8 passwords were cracked on average, compared to only 25.8% comprehensiveSubset passwords. The two test sets diverge as early as one billion guesses (6.8% to 0.5%).

Ignoring comprehensiveSubset passwords that were created under the basic16 condition allows us to analyze 171 passwords, all created under less strict conditions. Only 25.2% of these passwords are cracked on average, suggesting that subsets drawn exclusively from less strict conditions are more difficult to guess than passwords created under stricter requirements.

To understand this result more deeply, we examined the distribution of structures in the two test sets. There are 618 structures in the 1000-password comprehensive8 set, compared to 913 for comprehensiveSubset (normalized). Fifty-two percent of comprehensive8 passwords have unique structures, compared to 85% for comprehensiveSubset. This distribution of structures explains why comprehensive8 is significantly easier to guess.

We do not know why the two samples are different, although we suspect it may be related to the comprehensiveSubset subset isolating those users who make the most complex passwords. Regardless of the reason for this difference, however, researchers seeking to compare password policies should be aware that such subsets may not be representative.

5.4 Guessability and entropy

Historically, Shannon entropy (computed or estimated using various methods) has provided a convenient single statistic to summarize password strength. It remains unclear, however, how well entropy reflects the guess resistance of a password set. While information entropy does provide a theoretical lower bound on the guessability of a set of passwords [28], in practice a system administrator may be more concerned about how many passwords can be cracked in a given number of guesses than about the average guessability across the population. Although there is no mathematical relationship between entropy and this definition of guess resistance, we examine the possibility that the two are correlated in practice. To do this, we consider two independent measures of entropy, as defined in Section 4.2: an empirically calculated estimate and a theoretical NIST estimate. For both measures, we find that entropy estimates roughly indicate which composition policies provide more guess resistance than others, but provide no useful information about the magnitude of these differences.

Empirically estimated entropy. We ranked our password conditions based on the proportion of passwords cracked in our most complete experiment (E) at one trillion guesses, and compared this to the rank of conditions based on empirically estimated entropy. We found these rankings, shown in Figure 9, to be significantly correlated (Kendall’s $\tau = 0.71$, Holm-corrected $p = 0.042$). However, looking at the proportion of passwords cracked at a million or a billion guesses, the correlation in rankings is no longer significant (Holm-corrected $p = 0.275, 0.062$). The same pattern of significance, correlation at one trillion guesses but not at one billion or one million, was found in our largest

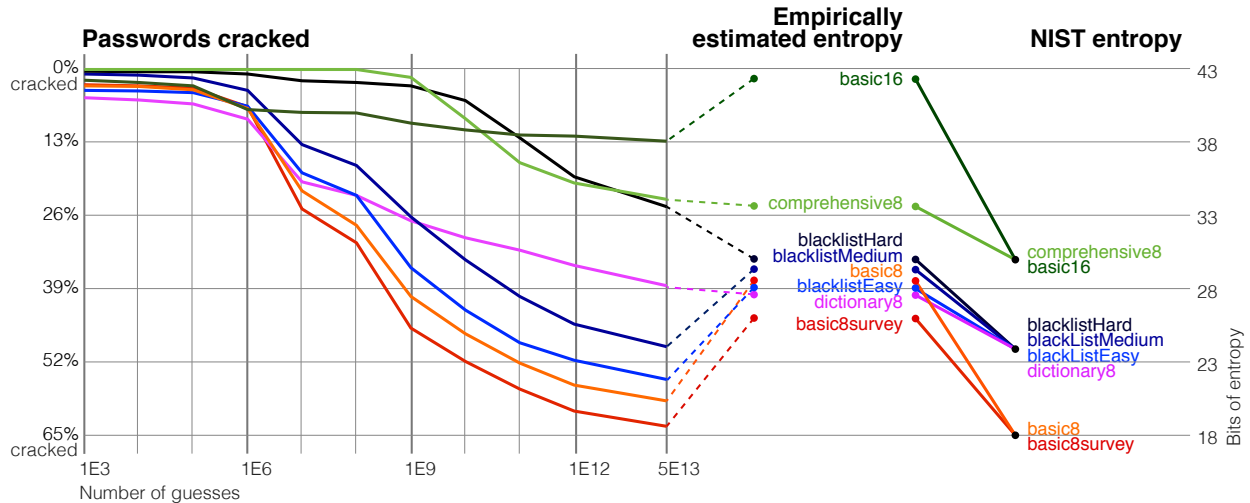


Figure 9: Relationship among the resistance of our collected password sets to heuristic cracking (experiment E); empirical entropy estimates we calculate from those sets; and NIST entropy estimates for our password conditions.

public-data experiment (P4). These results indicate that entropy might be useful when considering an adversary who can make a large number of guesses, but is not useful when considering a smaller number of guesses.

Further, empirically estimated entropy was unable to predict correctly the ranking of dictionary8, even when considering a large number of guesses. This condition displayed greater resistance to guessability than basic8, yet its empirically estimated entropy was lower. This might indicate a flaw in how entropy was estimated, a flaw in the guessing algorithm, or an innate shortcoming of the use of entropy to predict guessability. Since entropy can only lower-bound the guessability of passwords, it is possible for the frequency distribution of dictionary8 to have low entropy but high guess resistance. If this is the case, Verheul theorized that such a distribution would be optimal for password policy [43].

NIST entropy. Computing the NIST entropy of our password conditions produces three equivalence classes, as shown in Figure 9. These arise because NIST entropy is not granular enough to capture all differences between our conditions. First, NIST entropy does not take into account the size of a dictionary or its implementation. All five of our dictionary and blacklist conditions meet the NIST requirement of a dictionary with at least 50,000 words [8]. Implementation details, such as case-insensitive blacklist checking or the removal of non-alphabetic characters before a dictionary check, are not considered in the entropy score. Our results show that these details lead to password policies with very different levels of password strength and should be considered in a future heuristic.

Second, the NIST entropy scores for basic16 and comprehensive8 are the same, even though basic16 appears to be much more resistant to powerful guessing attacks. This may suggest that future heuristics should assign greater value to length than does the NIST heuristic.

Perhaps surprisingly, the equivalence classes given by NIST entropy are ordered correctly based on our results for guessability after 50 trillion guesses. Though its lack of granularity fails to capture differences between similar password conditions, NIST entropy seems to succeed at its stated purpose of providing a “rough rule of thumb” [8].

We stress that although both measures of entropy provide a rough ordering among policies, they do not always correctly classify guessability (see for example dictionary8), and they do not effectively measure how much additional guess resistance one policy provides as compared to another. These results suggest that a “rough rule of thumb” may be the limit of entropy’s usefulness as a metric.

6 Discussion

We next discuss a number of important issues regarding ethics, ecological validity, and the limitations of our methodology.

Ethical considerations. Most of our results rely on passwords we have collected during a user study (approved by our institution’s IRB). However, we also use the RockYou and MySpace password lists. Although these passwords have collectively been used by a number of scientific works that study passwords (e.g., [4, 13, 46, 47]), this nevertheless creates an ethical conundrum: Should our research use passwords acquired illicitly? Since this data has already been

made public and is easily available, using it in our research does not increase the harm to the victims. We use these passwords only to train and test guessing algorithms, and not in relationship with any usernames or other login information. Furthermore, as attackers are likely to use these password sets as training sets or cracking dictionaries, our use of them to evaluate password strength implies our results are more likely to be of practical relevance to security administrators.

Ecological validity. As with any user study, the ecological validity of our approach is important to the generalizability of our results. First, it is important to understand the results in the context of our participant sample. As we describe in Sections 1 and 3.4, our sample of Mechanical Turk participants is somewhat younger and more educated than the general population, but more diverse than typical small-sample password studies.

A second factor inviting consideration is that the passwords we collected did not protect high-value accounts. As we describe in Section 1, this is a longstanding limitation of password research. To gain insight into the extent to which our participants behaved as they would in non-study conditions, we tested two password-creation scenarios (Section 3.3): one was taking a survey, designed to observe user behavior with passwords for short-term, low-value accounts; and one was a simulated change to a longer-term, higher-value email account. Our users provided stronger passwords (measured by guessability and entropy) in the email scenario, a result consistent with users picking better passwords to protect a (hypothetical) high-value e-mail account than a low-value survey account. All our conditions except basic8 used the email scenario.

In our study, as in the real world, some users wrote down or otherwise stored their passwords. We asked participants who returned for the second half of the study whether or not they stored the password they had created (after reassuring them they would get paid either way), and we also instrumented the password-entry form to detect copy-paste and browser auto-fill behavior. We detected about 6% of participants using these methods of storage, while overall about one third admitted storing their passwords. Participants in comprehensive8 stored their passwords significantly more often than those in the other conditions (PHFET, $p < 0.05$).

We designed our study to minimize the impact of sampling and account-value limitations. All our findings result from comparisons *between* conditions. *caused by* the ways in which conditions differ (e.g., using a different technique to choose longer passwords than shorter ones) would be correctly captured and appropriately reflected in the results. Thus, we believe it is likely that our findings hold in general, at least for some classes of passwords and some classes of users.

Other limitations. We tested all sets of passwords with a number of password-guessing tools; the one we focus on (the Weir algorithm) always performed best. There may exist algorithms or training sets that would be more effective at guessing passwords than anything we tested. While this might affect some of our conclusions, we believe that most of them are robust, partly because many of our results are supported by multiple experiments and metrics.

In this work, we focused on automated offline password-guessing attacks. There are many other real-life threats to password security, such as phishing and shoulder surfing. Our analysis of password strength does not account for these. The password-composition policies we tested may induce different behaviors, e.g., writing down or forgetting passwords or using password managers, that affect password security. Although such effects have previously been studied for a subset of the policies in this study [25], space constraints dictate that a comprehensive investigation is beyond the scope of this paper.

7 Conclusion

Although the number and complexity of password-composition requirements imposed by systems administrators at a wide range of organizations have been steadily increasing, the actual value added by these requirements is poorly understood. In this work, we take a substantial step forward in understanding not only these requirements themselves, but also the process of evaluating them.

We introduced a new, efficient technique for evaluating password strength that can be implemented for a variety of password-guessing algorithms and tuned using a variety of training sets to gain insight into the comparative guess resistance of different sets of passwords. Using this technique, we were able to perform a more comprehensive password analysis than had previously been possible.

We found several notable results about the comparative strength of different composition policies. Although NIST considers basic16 and comprehensive8 equivalent, we found that basic16 is superior for large numbers of guesses. Combined with a prior result that basic16 is also easier for users, this suggests that basic16 is the better policy choice [25]. We also found that the effectiveness of a dictionary check depends heavily on the choice of dictionary; in particular, a large blacklist created using state-of-the-art password-guessing techniques is much more effective than

a standard dictionary at preventing users from choosing easily guessed passwords. Our findings highlight several interesting points in the password-policy space and suggest some directions for further research to more fully detail a complete set of tradeoffs among composition-policy requirements.

Our results also reveal important information about conducting guess-resistance analysis. Effective attacks on passwords created under complex or rare-in-practice composition policies require access to abundant, closely matched training data. In addition, this type of password set cannot be characterized correctly simply by selecting a subset of conforming passwords from a larger corpus; such a subset is unlikely to be representative of passwords created under the policy in question. Finally, we report that Shannon entropy, though a convenient single-statistic metric of password strength, provides only a rough correlation with guess resistance and is unable to correctly predict quantitative differences in guessability among password sets.

8 Acknowledgments

We thank Mitch Franzos of the Carnegie Mellon University Parallel Data Laboratory. This research was supported by NSF grants DGE-0903659, CCF-0424422, and CNS-111676, by CyLab at Carnegie Mellon under grants DAAD19-02-1-0389 and W911NF-09-1-0273 from the Army Research Office, and by a gift from Microsoft Research.

References

- [1] ADAMS, A., SASSE, M. A., AND LUNT, P. Making passwords secure and usable. In *HCI 97* (1997).
- [2] BISHOP, M., AND KLEIN, D. V. Improving system security via proactive password checking. *Computers & Security* 14, 3 (1995), 233–249.
- [3] BONNEAU, J. The Gawker hack: how a million passwords were lost, Dec. 2010. <http://www.lightbluetouchpaper.org/2010/12/15/the-gawker-hack-how-a-million-passwords-were-lost/>.
- [4] BONNEAU, J., JUST, M., AND MATTHEWS, G. What’s in a name? evaluating statistical attacks on personal knowledge questions. In *Proc. Financial Crypto. 2010* (Tenerife, Spain, Jan. 2010), pp. 98–113.
- [5] BRIGHT, P. Anonymous speaks: The inside story of the HBGary hack. <http://arstechnica.com/tech-policy/news/2011/02/anonymous-speaks-the-inside-story-of-the-hbgary-hack.ars>, February 2011.
- [6] BRIGHT, P. “Military Meltdown Monday”: 90K military usernames, hashes released. <http://arstechnica.com/tech-policy/news/2011/07/military-meltdown-monday-90k-military-usernames-hashes-released.ars>, July 2011.
- [7] BUHRMESTER, M., KWANG, T., AND GOSLING, S. D. Amazon’s Mechanical Turk: A new source of inexpensive, yet high-quality, data? *Perspectives on Psychological Science* 6, 1 (2011), 3–5.
- [8] BURR, W. E., DODSON, D. F., AND POLK, W. T. Electronic authentication guideline. Tech. rep., NIST, 2006.
- [9] CHIASSON, S., FORGET, A., STOBERT, E., VAN OORSCHOT, P. C., AND BIDDLE, R. Multiple password interference in text passwords and click-based graphical passwords. In *Proceedings of the 16th ACM conference on Computer and communications security* (New York, NY, USA, 2009), CCS ’09, ACM, pp. 500–511.
- [10] CONSTANTIN, L. Sony Stresses that PSN Passwords Were Hashed. <http://news.softpedia.com/news/Sony-Stresses-PSN-Passwords-Were-Hashed-198218.shtml>, May 2011.
- [11] DAVIS, D., MONROSE, F., AND REITER, M. K. On user choice in graphical password schemes. In *Proceedings of the 13th conference on USENIX Security Symposium - Volume 13* (Berkeley, CA, USA, 2004), SSYM’04, USENIX Association, pp. 11–11.
- [12] DEAN, J., AND GHEMAWAT, S. MapReduce: Simplified Data Processing on Large Clusters. In *Symp. on Operating System Design and Implementation (OSDI)* (San Francisco, CA, Dec 2004).

- [13] DELL'AMICO, M., MICHIARDI, P., AND ROUDIER, Y. Password strength: An empirical analysis. In *Proc. INFOCOM 2010* (San Diego, CA, Mar. 2010), pp. 983–991.
- [14] DESIGNER, S. John the Ripper. <http://www.openwall.com/john/>, 1996-2010.
- [15] DOWNS, J. S., HOLBROOK, M. B., SHENG, S., AND CRANOR, L. F. Are your participants gaming the system?: screening mechanical turk workers. In *Proceedings of the 28th international conference on Human factors in computing systems* (New York, NY, USA, 2010), CHI '10, ACM, pp. 2399–2402.
- [16] FLORÊNCIO, D., AND HERLEY, C. A large-scale study of web password habits. In *Proc. WWW'07* (2007).
- [17] FLORÊNCIO, D., AND HERLEY, C. Where do security policies come from? In *Proc. SOUPS '10* (2010).
- [18] GAW, S., AND FELTEN, E. W. Password management strategies for online accounts. In *Proceedings of the second symposium on Usable privacy and security* (New York, NY, USA, 2006), SOUPS '06, ACM, pp. 44–55.
- [19] HART, D. Attitudes and practices of students towards password security. *Journal of Computing Sciences in Colleges* 23, 5 (2008), 169–174.
- [20] INGLESANT, P., AND SASSE, M. A. The true cost of unusable password policies: password use in the wild. In *Proc. ACM CHI'10* (2010), pp. 383–392.
- [21] IPEIROTIS, P. G. Demographics of mechanical turk. Tech. Rep. CeDER-10-01, New York University, March 2010.
- [22] KELLEY, P. G. Conducting usable privacy and security studies with Amazon's Mechanical Turk. In *Proceedings of the USER Workshop at the Symposium on Usable Privacy and Security (SOUPS)* (2010), USER 2010.
- [23] KITTUR, A., CHI, E. H., AND SUH, B. Crowdsourcing user studies with mechanical turk. In *Proceeding of the twenty-sixth annual SIGCHI conference on Human factors in computing systems* (New York, NY, USA, 2008), CHI '08, ACM, pp. 453–456.
- [24] KOMANDURI, S., AND HUTCHINGS, D. R. Order and entropy in picture passwords. In *Graphics Interface* (2008), pp. 115–122.
- [25] KOMANDURI, S., SHAY, R., KELLEY, P. G., MAZUREK, M. L., BAUER, L., CHRISTIN, N., CRANOR, L. F., AND EGELMAN, S. Of passwords and people: Measuring the effect of password-composition policies. In *CHI 2011: Conference on Human Factors in Computing Systems* (May 2011).
- [26] KUO, C., ROMANOSKY, S., AND CRANOR, L. F. Human selection of mnemonic phrase-based passwords. In *Symposium on Usable Privacy and Security* (2006), pp. 67–78.
- [27] MARECHAL, S. Advances in password cracking. *Journal in Computer Virology* 4, 1 (2008), 73–81.
- [28] MASSEY, J. L. Guessing and entropy. In *Proc. IEEE Int. Symp. Info. Theory* (1994), p. 204.
- [29] NARAYANAN, A., AND SHMATIKOV, V. Fast dictionary attacks on passwords using time-space tradeoff. In *CCS '05: Proceedings of the 12th ACM conference on Computer and communications security* (New York, NY, USA, 2005), ACM, pp. 364–372.
- [30] PROCTOR, R. W., LIEN, M.-C., VU, K.-P. L., SCHULTZ, E. E., AND SALVENDY, G. Improving computer security for authentication of users: Influence of proactive password restrictions. *Behavior Res. Methods, Instruments, & Computers* 34, 2 (2002), 163–169.
- [31] ROSS, J., IRANI, L., SILBERMAN, M. S., ZALDIVAR, A., AND TOMLINSON, B. Who are the crowdworkers?: shifting demographics in mechanical turk. In *Proceedings of the 28th of the international conference extended abstracts on Human factors in computing systems* (New York, NY, USA, 2010), CHI EA '10, ACM, pp. 2863–2872.
- [32] SCHECHTER, S., HERLEY, C., AND MITZENMACHER, M. Popularity is everything: A new approach to protecting passwords from statistical-guessing attacks. In *Proc. HotSec'10* (2010).

- [33] SCHNEIER, B. Myspace passwords aren't so dumb. <http://www.wired.com/politics/security/commentary/securitymatters/2006/12/72300>, December 2006, retrieved November 2010.
- [34] SHANNON, C. E. A mathematical theory of communication. *Bell Syst. Tech. J.* 27 (1949), 379–423,623–656.
- [35] SHANNON, C. E. Prediction and entropy of printed english. *Bell Systems Tech. J.* 30 (1951), 50–64.
- [36] SHAY, R., AND BERTINO, E. A comprehensive simulation tool for the analysis of password policies. *Int. J. Info. Sec.* 8, 4 (2009), 275–289.
- [37] SHAY, R., BHARGAV-SPANTZEL, A., AND BERTINO, E. Password policy simulation and analysis. In *ACM workshop on Digital identity management* (2007), pp. 1–10.
- [38] SHAY, R., KOMANDURI, S., KELLEY, P., LEON, P., MAZUREK, M., BAUER, L., CHRISTIN, N., AND CRANOR, L. Encountering stronger password requirements: user attitudes and behaviors. In *Proc. SOUPS'10* (2010).
- [39] SPAFFORD, E. H. OPUS: Preventing weak password choices. *Computers & Security* 11, 3 (1992), 273–278.
- [40] STANTON, J. M., STAM, K. R., MASTRANGELO, P., AND JOLTON, J. Analysis of end user security behaviors. *Comp. & Security* 24, 2 (2005), 124 – 133.
- [41] TOOMIM, M., KRIPLEAN, T., PÖRTNER, C., AND LANDAY, J. Utility of human-computer interactions: toward a science of preference measurement. In *Proceedings of the 2011 annual conference on Human factors in computing systems* (New York, NY, USA, 2011), CHI '11, ACM, pp. 2275–2284.
- [42] VANCE, A. If your password is 123456, just make it hackme. *New York Times*, <http://www.nytimes.com/2010/01/21/technology/21password.html>, January 2010, retrieved September 2010.
- [43] VERHEUL, E. Selecting secure passwords. In *Topics in Cryptology – CT-RSA 2007*, M. Abe, Ed., vol. 4377 of *Lecture Notes in Computer Science*. Springer Berlin / Heidelberg, 2006, pp. 49–66. 10.1007/11967668_4.
- [44] VU, K.-P. L., PROCTOR, R. W., BHARGAV-SPANTZEL, A., TAI, B.-L. B., AND COOK, J. Improving password security and memorability to protect personal and organizational information. *Int. J. of Human-Comp. Studies* 65, 8 (2007), 744–757.
- [45] WEIR, C. M. *Using Probabilistic Techniques To Aid In Password Cracking Attacks*. PhD thesis, 2010.
- [46] WEIR, M., AGGARWAL, S., COLLINS, M., AND STERN, H. Testing metrics for password creation policies by attacking large sets of revealed passwords. In *CCS '10: Proceedings of the 17th ACM conference on Computer and communications security* (New York, NY, USA, 2010), ACM, pp. 162–175.
- [47] WEIR, M., AGGARWAL, S., DE MEDEIROS, B., AND GLODEK, B. Password cracking using probabilistic context-free grammars. In *2009 30th IEEE Symposium on Security and Privacy* (2009), IEEE, pp. 391–405.
- [48] WHITE, T. *Hadoop: The Definitive Guide*, 2nd ed. O'Reilly, September 2010.
- [49] WIMBERLY, H., AND LIEBROCK, L. M. Using fingerprint authentication to reduce system security: An empirical study. In *Security and Privacy (SP), 2011 IEEE Symposium on* (may 2011), pp. 32 –46.
- [50] ZHANG, Y., MONROSE, F., AND REITER, M. K. The security of modern password expiration: an algorithmic framework and empirical analysis. In *CCS '10: Proceedings of the 17th ACM conference on Computer and communications security* (New York, NY, USA, 2010), ACM, pp. 176–186.
- [51] ZVIRAN, M., AND HAGA, W. J. Password security: an empirical study. *J. Mgt. Info. Sys.* 15, 4 (1999), 161–185.

A Calculator Experiments

Here we detail the complete training and test data used in each of our Weir-algorithm experiments. The first column gives the experiment number. The next three columns list the three types of training data used to create a Weir-calculator experiment. The *structures* column describes the wordlist(s) used to generate the set of character-type structures that define the Weir algorithm’s search space. The *digits and symbols* column lists the wordlist(s) that determine the probabilities with which combinations of digits and symbols can be filled into those structures. The *strings* column shows which wordlists provide the probabilities with which alphabetic strings are filled into structures. In most cases, we train strings on as much data as possible, while restricting structure and digit/symbol training to those wordlists that contain a quality sample of multi-character-class passwords. In the final column, we describe the set(s) of passwords that we attempted to guess in a given experiment.

We also list the complete training and test data used in each of our BFM experiments. The experiment number and test set columns are the same as in the Weir subtable. Training for the BFM calculator, however, is considerably simpler, using only one combined wordlist per experiment; these lists are detailed in the *training set* column.

Abbreviations for all the training and test sets we use are defined in the key below the tables.

Weir experiment descriptions

Name	Training sets			Testing Set
	Structures	Digits and symbols	Strings	
Trained from public password data				
P1	MS8	MS	MS	1000-All
P2	MS8	MS	MS, W2, I	1000-All
P3	MS8	MS, RY	MS, W2, I, RY	1000-All
P3-C8	MSC	MS, RY	MS, W2, I, RY	1000-C8
P3-B16	MS16	MS, RY	MS, W2, I, RY	1000-B16
P4	MS8, OW8	MS, RY, OW	MS, W2, I, RY, OW	1000-All
P4-B16	MS16, OW16	MS, RY, OW	MS, W2, I, RY, OW	1000-B16
Trained on half of our dataset, weighted to 1/10th, equal-size, or 10x the cumulative size of the public data				
X1/10	MS8, 500-All	MS, RY, 500-All	MS, W2, I, RY, 500-All	500-All
X1	MS8, 500-All	MS, RY, 500-All	MS, W2, I, RY, 500-All	500-All
X10	MS8, 500-All	MS, RY, 500-All	MS, W2, I, RY, 500-All	500-All
Everything				
E	MS8, OW8, 500-All	MS, RY, OW, 500-All	MS, W2, I, RY, OW, 500-All	500-All
Testing password subsets that meet comprehensive8 requirements				
S0a	MSC, OWC	MS, OW	MS, W2, I, OW	1000-C8, 206-C8S
S0b	MSC, OWC, 2000-C8	MS, OW, 2000-C8	MS, W2, I, OW, 2000-C8	1000-C8, 206-C8S
S1	MSC, OWC, RYCD	MS, OW, RY	MS, W2, I, OW, RY	1000-C8, 206-C8S
S2	MSC, OWC, 2000-C8, RYCD	MS, OW, 2000-C8, RY	MS, W2, I, OW, 2000C8, RY	1000-C8, 206-C8S
Split ratio testing on basic8				
B8a	MS8, OW8, 500-B8	MS, RY, OW, 500-B8	MS, W2, I, RY, OW, 500-B8	2500-B8
B8b	MS8, OW8, 1000-B8	MS, RY, OW, 1000-B8	MS, W2, I, RY, OW, 1000-B8	2000-B8
B8c	MS8, OW8, 1500-B8	MS, RY, OW, 1500-B8	MS, W2, I, RY, OW, 1500-B8	1500-B8
B8d	MS8, OW8, 2000-B8	MS, RY, OW, 2000-B8	MS, W2, I, RY, OW, 2000-B8	1000-B8
B8e	MS8, OW8, 2500-B8	MS, RY, OW, 2500-B8	MS, W2, I, RY, OW, 2500-B8	500-B8
Split ratio testing on comprehensive8				
C8test1/10	MSC, 500-C8	MS, RY, 500-C8	MS, W2, I, RY, 500-C8	2500-C8
C8test1	MSC, 500-C8	MS, RY, 500-C8	MS, W2, I, RY, 500-C8	2500-C8
C8a	MSC, OWC, 500-C8	MS, RY, OW, 500-C8	MS, W2, I, RY, OW, 500-C8	2500-C8
C8b	MSC, OWC, 1000-C8	MS, RY, OW, 1000-C8	MS, W2, I, RY, OW, 1000-C8	2000-C8
C8c	MSC, OWC, 1500-C8	MS, RY, OW, 1500-C8	MS, W2, I, RY, OW, 1500-C8	1500-C8
C8d	MSC, OWC, 2000-C8	MS, RY, OW, 2000-C8	MS, W2, I, RY, OW, 2000-C8	1000-C8
C8e	MSC, OWC, 2500-C8	MS, RY, OW, 2500-C8	MS, W2, I, RY, OW, 2500-C8	500-C8

BFM experiment descriptions

Name	Training set	Test set
B1	RY, MS, I	1000-All
B2	RY, MS, I, 500-All	500-All
B3	RY, MS, I, 2000-B8	1000-B8
B4	RY, MS, I, 2000-C8	1000-C8

Key to password sets

RY	RockYou list	I	inflection list
RYCD	RY, filtered w/ all reqs. of C8	W2	simple Unix dictionary
MS	MySpace list	OW	paid Openwall dictionary
MS8	MS, filtered w/ min length of 8	OW8	OW, filtered w/ min length of 8
MS16	MS, filtered w/ min length of 16	OW16	OW, filtered w/ min length of 16
MSC	MS, filtered w/ min length of 8 and character class reqs. of C8	OWC	OW, filtered w/ min length 8 and character class reqs. of C8
<i>n</i>-All	<i>n</i> passwords from each of our conditions	<i>n</i>-B8	<i>n</i> basic8 passwords
<i>n</i>-B16	<i>n</i> basic16 passwords	<i>n</i>-C8	<i>n</i> comprehensive8 passwords
<i>n</i>-C8S	<i>n</i> comprehensiveSubset passwords	<i>n</i>-RYCD	<i>n</i> RYCD passwords