

Please Continue to Hold

An empirical study on user tolerance of security delays

Serge Egelman
Brown University
egelman@cs.brown.edu

David Molnar
Microsoft Research
dmolnar@microsoft.com

Nicolas Christin
Carnegie Mellon University
nicolasc@andrew.cmu.edu

Alessandro Acquisti
Carnegie Mellon University
acquisti@andrew.cmu.edu

Cormac Herley
Microsoft Research
cormac@microsoft.com

Shriram Krishnamurthi
Brown University
sk@cs.brown.edu

ABSTRACT

We present the results of an experiment examining the extent to which individuals will tolerate delays when told that such delays are for security purposes. In our experiment, we asked 400 Amazon Mechanical Turk users to count the total number of times a certain term was repeated in a multi-page document. The task was designed to be conducive to cheating. We assigned subjects to four between-subjects conditions: one of these offered a concrete security reason (virus-scanning) for the delay, another offered only a vague security explanation, another did not mention security at all, and a control condition without a delay.

We found that subjects were significantly more likely to cheat or abandon the task when provided no explanation or a vague security explanation for the delay. However, when subjects were provided more explanation about the threat model and the protection ensured by the delay, they were not more likely to cheat than subjects in the control condition who faced no such delay. Our results thus contribute to the nascent literature on soft paternalistic solutions to security and privacy problems by suggesting that, when security mitigations cannot be made “free” for users, designers may incentivize compliant users’ behavior by intentionally drawing attention to the mitigation itself.

1. INTRODUCTION

The computer security community has proposed various approaches to thwarting security breaches. Such *security mitigations* can be divided into three different categories: mitigations that are invisible to the user, mitigations that do something noticeable on the user’s behalf, and mitigations that suggest the user take a particular action. The latter category, in particular, can come at a substantial cost for the user. If this cost is too large, users may choose to forgo the security mitigation, with detrimental individual and collective effects: negative externalities may be created when users fall victim to security breaches [13], and revenues for a firm producing more secure products may be lost if users choose less burdensome (even though less secure) software.

In this paper we explore the cost of security mitigations from the perspective of the end user, and how manipulating – and offering explanations for – those costs may impact users’ acceptance of security mitigations. We are specifically interested in quantifying the inconveniences users will accept in the name of security. Previous work on security mitigations investigated costs [5, 6, 8, 12, 16, 21, 24, 25, 32], but

focused on application compatibility and the speed impact of the mitigation on a set of benchmarks. In contrast, we attempt to measure what makes costs of mitigations acceptable to users. We posit that not only should acceptability of security mitigations be evaluated directly through user studies, but that non-normative mechanisms should be devised to increase the acceptability of those mitigations. By “non-normative,” we refer to mechanisms that do not affect the technical performance of the security mitigation, but may influence the way users react to it: for instance, providing explanations with varying degrees of detail, making certain types of information more or less salient, or artificially manipulating speed and delays in the product’s performance.

As a first step in this line of inquiry, we designed an experiment in which we asked subjects to count the number of times a specific word appeared in a PDF document using a custom Flash-based viewer. We purposefully designed this task to be conducive to *cheating*: the user could submit a response without actually reading the entire document. The rate of cheating with a “mitigation” in place compared to “no mitigation” then gives us a quantitative measure of the acceptability of a security mitigation. That is, if it takes an unacceptable amount of time to complete the task in one condition, we would expect to see a disproportionate amount of cheating in that condition.

We controlled for two common user-facing aspects of a security mitigation: the presence of a delay and the presence of a notice explaining the reason for the delay. We designed our study as a four condition between-subjects experiment: one control condition, with no delay and no notice; one “loading” condition, with a delay but no notice; one condition with a delay and a vague notice; and one condition that incorporated a delay along with a detailed explanation of the threat model and its mitigation. Based on results from the behavioral economics and social psychology literature, we hypothesized that subjects would be more tolerant of the delay when told that the delay was for security purposes, and when primed with more detailed explanations of the threat model and its mitigation. Thus, we expected to see a disproportionate amount of cheating in the condition with the delay and no notice.

We conducted our experiment with 400 subjects recruited from Amazon’s Mechanical Turk. Mechanical Turk is a service for advertising “tasks” that can be completed by human workers for a set price per task. We found that subjects were significantly more likely to cheat or abandon the task

when provided no explanation or a vague security explanation for the delay. However, when subjects were provided more explanation about the threat model and the protection ensured by the delay, they were not more likely to cheat than subjects in the control condition who faced no such delay.

Our findings suggest that if a delay is necessary due to a security mitigation, then the mitigation may be more acceptable to users if they are told of the threat model and how the mitigation protects them. This advice may sound counterintuitive, given that most security mitigations intentionally do not make themselves visible to the user unless an attack occurs, and perhaps not even then. For example, Windows programs do not usually employ pop-ups to inform users that address space randomization is being used or that stack canaries are being inserted. On the other hand, many anti-virus programs do explicitly warn the user when a scan is in progress, which may incur significant delays to user operations. There is therefore a tension between security which is “invisible,” and measures made explicit “for security reasons.” In this paper, we attempt to better characterize these tensions by providing empirical evidence. Our work shows that user studies of security mitigations shed important light on the acceptability of mitigations in contexts not well served by previous approaches. In this regard, our results contribute to the nascent literature on soft paternalistic solutions to security and privacy problems.

Our approach side-steps the question of “is $X\%$ overhead on this benchmark a lot or a little?” by placing the user-facing aspects of the mitigation directly in the context in which they are experienced by users. Because the technical aspects of the actual mitigation can be abstracted away, these types of user studies can be generalized to yield actionable findings for multiple types of security mitigations. Therefore, we believe user studies are an important addition to traditional benchmarking and application compatibility analysis for evaluating security mitigations. Our experience with Mechanical Turk shows that these studies can be carried out at modest cost even with hundreds of users. We hope this will encourage others to perform such studies as part of the process for evaluating future security mitigations.

2. BACKGROUND

Our work responds to and is informed by traditions from computer security, behavioral economics, human-computer interaction/computer-supported collaborative work, and psychology. We now discuss background from each of these communities in detail.

Computer Security. *Security mitigations* are features of an application or operating system that make it more difficult for an adversary to take control of a victim’s computer, even when the victim’s software has a bug such as a buffer overflow. The value of a security mitigation is that it trades speed and application compatibility for increased attack difficulty, *without requiring the defender to have detailed, specific knowledge of the attack in advance*. This tradeoff is attractive because finding all bugs in computer software or enumerating specific attacks ahead of time is incredibly difficult.

The computer security community has a long list of proposed mitigations stretching over more than fifteen years. Classic examples include stack canaries [6,12], address space layout randomization [24], automatic bounds checking [5, 16], and non-executable memory [21,24]. More recent exam-

ples include software changes such as Nozzle [25], or efficient software fault isolation for x86 and x64 architectures [20,32], and hardware changes such as those found in Raksha [8] or SmashGuard [23]. Anti-virus software can also be viewed in this category, and it has been undeniably successful commercially.

In the evaluation section of every proposal of which we are aware, the cost of a mitigation is evaluated along two axes: the loss in speed on benchmarks and the impact on application compatibility. Application compatibility is an important consideration, but one which we do not address in this work. Another increasingly important phenomenon is that adversaries can develop reliable methods for bypassing mitigations [10]; once such bypasses are found, users pay the cost of the mitigation but receive no benefit against sophisticated adversaries.

Even so, a major part of discussions on adoption from the earliest mitigations to the present centers on whether a specific speed impact is “too much,” as measured on a set of benchmarks. The key problem we address is that measuring speed on benchmarks is merely a proxy for measuring *user acceptance* of a mitigation. Clearly, if there is no speed impact from adopting a mitigation, the experience of the user with the mitigation is indistinguishable from the original experience in the common case where no attack is present. Therefore the mitigation will be acceptable.

Unfortunately, in most cases security mitigations cannot be made “free.” The computer security community has then historically proceeded to the difficult question: “is the speed impact of the mitigation on these benchmarks acceptable?” Complicating the question is the fact that “acceptance” means different things in different scenarios. For example, an additional 50 milliseconds to load a web page may lead to a significant loss in revenue for a web site. An additional hour added to a batch job, which normally takes a year, may not be noticed.

Behavioral Economics, Usability, and Soft Paternalism. Our contribution can be related to the nascent literature on the application of soft paternalistic approaches [19,30] to privacy [2,31] and security [4] problems. In recent years, there has been growing interest in understanding the psychological motives, as well as the possible cognitive and behavioral biases, that affect privacy [1] and security [27] decision making. In parallel, the computer science community has started investigating how to make privacy and security systems more usable [7]. These streams of research converge when lessons derived from the behavioral economics, decision research, or psychological literatures are incorporated into the design of systems that take into account systematic biases that affect our decision making in information security. These approaches do not merely aim at making systems more usable, but to actually anticipate known and costly biases – and sometimes even exploit those biases in manners that nudge users towards certain choices, without limiting their freedom [19,30]: from providing salient information [31] to creating interactive audited dialogs [4] with the end users.

Psychology. Milgram explored the role of obedience to authority in his seminal 1963 experiment. He concluded that people are generally compliant with requests—however bizarre or nonsensical—when those requests come from people in positions of authority [22]. While work on obedience to computer security mitigations has some parallels, our work

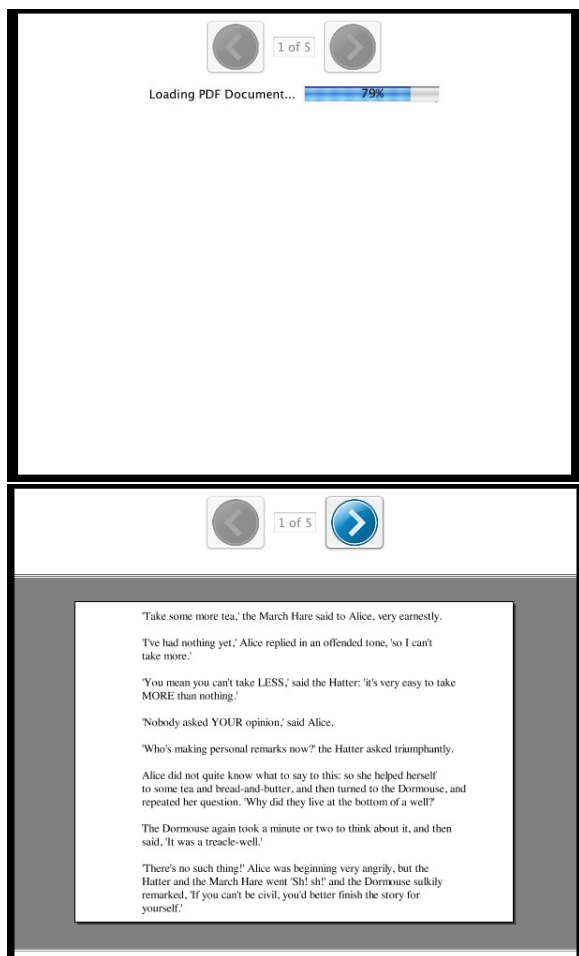


Figure 1: Screenshots of the SuperViewer interface. Subjects were able to change pages using the arrow buttons. The above image shows the loading screen that all subjects saw when first launching SuperViewer; subjects in the *control* condition never saw this screen again, while those in the experimental conditions repeatedly saw it on every page. The page number indicator was subsequently removed after our pilot studies.

differs in that our requests did not come from a human being in a position of authority. Our experiment tests a hypothesis related, in part, to the results of a famous experiment by Langer, Blank, and Chanowitz [18]. In their series of experiments, they showed that even “placebic” information provided in the form of explanation, or reason for a request, was sufficient to generate compliance with a request, even though the reason itself conveyed no actual information. In our experiments, we tested whether providing an explanation for the security delays in loading pages would, in fact, increase the likelihood that subjects would comply with the security mitigation.

3. METHODOLOGY

We conducted an experiment using Mechanical Turk to examine whether people would put up with an inconvenience if they were told it was for security purposes. Researchers

have recently begun using Mechanical Turk as a way to quickly perform large-scale human subjects experiments for very little cost [17]. In 2009, Ross et al. performed a series of surveys using Mechanical Turk and concluded that the demographics do not significantly differ from the population of U.S. Internet users [26]. In 2009, Jakobsson performed an experiment to study the quality of work produced by Mechanical Turk users. He commissioned a survey using Mechanical Turk as well as an identical one using an “established, independent survey company” and found no significant differences between the two participant pools [15].

We created a task wherein we told study subjects that they were beta testing a new web-based document viewer, *SuperViewer* (Figure 1). When first launching SuperViewer, all subjects saw a progress bar that took ten seconds to load before displaying the first page of the document. Those in the *control* condition never saw this progress bar again, while those in the three experimental conditions repeatedly saw this progress bar each time they viewed a new page of the document (i.e., they had to wait ten seconds each time they turned to a new page). The three experimental conditions differed based on the text used to explain the reason for the progress bar.

For the task itself, we asked subjects to read a document using SuperViewer. We told subjects that to receive payment, they must report the frequency that a particular word occurred in that document. Thus, they would have to read the entire document in order to accurately answer the question. SuperViewer features a very basic interface: two buttons for navigating forward and backward in the document, which forced subjects to view the document pages in order. Likewise, there was no “search” functionality, otherwise completing the task would have been trivial. One goal of this task was to make it appear indistinguishable from other non-research Mechanical Turk tasks (e.g., product categorization, image labeling, etc.). Thus, if subjects did not believe they were engaged in a research study for the public good, they may have been more inclined to “cheat.” We defined cheating as submitting a response to receive credit without using SuperViewer to read the entire document (e.g., reading part or none of the document). Our main interest was to examine whether subjects’ cheating varied based on the four between-group conditions we created:

- **Control** — SuperViewer was launched when subjects clicked a button. Immediately after launching, a progress bar was displayed for ten seconds with the label, “Loading.” After this ten second period, the first page of the document was displayed. Subjects could change pages in the document by clicking one of two arrow buttons. Thus, subjects were forced to view pages in order.
- **Loading** — This condition was identical to the *control* condition, with one exception: when advancing to a subsequent page after the first, an additional ten second progress bar—also labeled “loading”—was displayed before subjects were allowed to view the next page. Subjects only saw these progress bars once for each new page; subjects would not see a progress bar again when flipping to a previously viewed page. The purpose of this condition was to examine whether study subjects would tolerate an unknown delay or whether they would cheat by quitting the task early and reporting an incorrect word frequency.

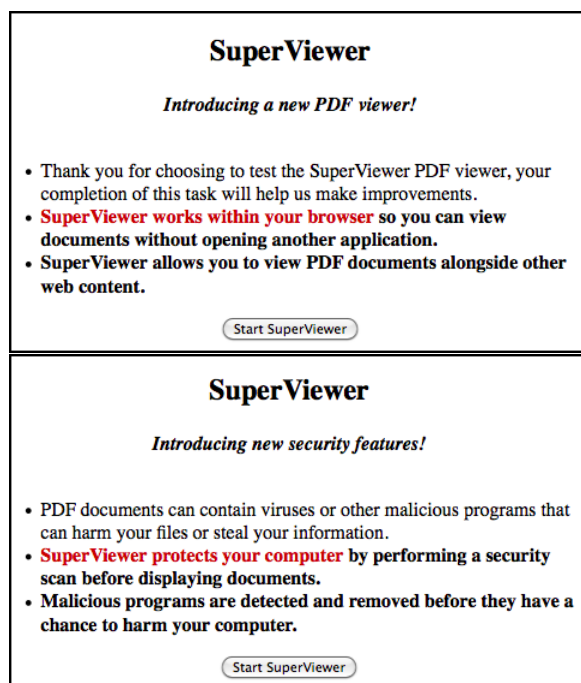


Figure 2: Screenshot of the instructions subjects saw before launching SuperViewer. The image above was seen by those in the *control*, *loading*, and *security* conditions. The image below was seen by those in the *primed* condition.

- **Security** — This condition was identical to the *loading* condition, with one exception: the label on the progress bar was changed from “Loading” to “Performing security scan.” The purpose of this condition was to examine whether study subjects would tolerate a delay when they were told it was for security purposes or if they would cheat by quitting the task early and reporting an incorrect word frequency.
- **Primed** — This condition was identical to the *security* condition, with one exception: prior to launching SuperViewer, subjects were informed of the danger of viruses embedded in online documents and that SuperViewer scans documents for their protection. The purpose of this condition was to examine whether subjects were any less likely to cheat if they understood why the “security scan” was being performed, or if simply performing an ambiguous security function was reason enough (as was the case in the *security* condition).

Prior to launching the SuperViewer applet, subjects were shown a page of information about the software (Figure 2). The main purpose of this page was to prime those in the *primed* condition to security concerns and to convince them that the software was protecting them against a legitimate threat model. In order to balance all of the conditions in terms of total workload required of our subjects—the total amount of text to read—we added a placebo text page for those in the other three conditions.

In order to determine whether our subjects cheated during the experiment, we recorded the number of unique document pages they viewed, the total number of pages viewed,

the time it took them between opening the document and submitting their response, and the numerical response that they submitted.

Upon completing the initial experiment, we invited subjects to complete a survey based on their experiences using SuperViewer in exchange for a bonus payment. The first page of the survey asked subjects about their overall opinions of SuperViewer and the factors that influenced those opinions: color, look and feel, ease of use, speed, and security features, each rated using a 5-point Likert scale. The second page of the survey contained questions about subjects’ risk perceptions, both when they used SuperViewer during the experiment, as well as when performing other activities on their computers (e.g., browsing the web, reading email, downloading files, etc.). The third page of the survey contained questions about what anti-virus software the subjects currently used, as well as the types of threats they believed said software guarded against. Finally, the fourth page of the survey featured demographic questions.

3.1 Pilot Studies

We decided to pilot our experiment using five pages from *Alice in Wonderland* as the document, and we used a version of SuperViewer written in Java. We offered to pay each participant \$0.05 to complete the task and we targeted 100 subjects, who were randomly assigned to the four conditions. Overall, we were underwhelmed at the rate of response to this task; it took us fifteen days to recruit 100 subjects. We decided that we needed to pay our subjects more, and so we created another 100 tasks, but this time paying \$0.11 per participant. This time it took us only six days to recruit 100 subjects.

While we had decided that we must pay subjects at least \$0.11 to complete this experiment in a timely manner, we discovered another potential caveat: several subjects had emailed us indicating that they could not load the viewer. In fact, while 200 people completed these pilot experiments, 355 others attempted to complete the task but were unsuccessful. Given that almost two thirds were unable to complete the task, we assumed that this was due to Java incompatibilities. Thus, we decided to rewrite SuperViewer in Flash.

We created a third pilot study to evaluate our Flash implementation. We recruited another 100 subjects and decided to pay them \$0.05, since we reasoned that with fewer technical incompatibilities from using Flash, we may receive an adequate participation rate with our original payment amount. Indeed, instead of taking fifteen days, using Flash allowed us to gather data from 100 subjects in just eight days. During the task itself, we asked subjects how many times the word “Hatter” occurred in the text, the correct answer being eight. Of our 100 subjects, only two clearly cheated, and each was in a different condition. We concluded that the task was too easy to perform, and therefore it would require an inordinate number of total subjects to get enough cheaters. Thus, we needed to make the task both longer and more frustrating; we changed the text from five pages of *Alice in Wonderland* to ten pages of *Ulysses*. The tenth page of the document only filled half a page, though we added a blank 11th page to indicate the end of the document. Because of this, we considered anyone who reached either page ten or eleven to have viewed the entire document.

Condition	N	Total Time (s)	Time Per Page (s)	Unique Pages	Total Pages	Cheaters
Control	119	458	50.9	9.93	19.76	15 (12.6%)
Loading	91	491	56.6	9.11	13.94	24 (26.4%)
Security	87	547	63.5	9.33	13.43	25 (28.7%)
Primed	103	546	56.9	10.01	15.10	18 (17.5%)

Table 1: Overview of our experimental data. This table shows the number of subjects in each of the four randomly-assigned conditions, the average time spent performing the task, the average reading time per page, the average number of unique pages viewed (out of a maximum of 11), the average number of total pages viewed, and finally the number of people who cheated by not reading the entire document before submitting a response.

The pilot versions of SuperViewer all displayed the current page number and the total number of pages in the document (Figure 1). By removing this status indicator, we reasoned that subjects will become more frustrated when they have no indication of when the task will end, and therefore, the observed effect size would be greater. Finally, we increased the payment to \$0.11 again, since we expected many more subjects to abandon the task without submitting any data, and therefore we needed a larger population.

4. ANALYSIS

A total of 400 Mechanical Turk users participated in our experiment between February 1st and 18th, 2010. These subjects were randomly assigned to the four conditions outlined in the previous section, with the total number of subjects in each condition presented in Table 1. As can be seen from the table, we observed several significant differences between the conditions. In this section we analyze these differences, both in terms of the number of people who cheated, as well as subjects’ task performance. Finally, we incorporate the data gleaned from our exit survey in order to support our conclusions.

4.1 Cheating

Subjects were required to read ten pages of *Ulysses* in order to answer the question, “how many times did the word *said* appear?” The correct answer was 23. We considered it cheating if a participant submitted an answer to this question without reading the entire ten pages. We hypothesized that subjects who had to wait for the progress bars to load before viewing subsequent pages would be significantly more likely to cheat than those in the *control* condition. We further hypothesized that subjects who were told that this delay was for “security purposes” would be less likely to cheat than those who were not given an explanation for the delay (i.e., those in the *loading* condition). Finally, we hypothesized that subjects who were given details of the threat model and the mitigation (i.e., the *primed* condition) would be just as likely to cheat as those who did not receive this information (i.e., the *security* condition).

Overall, we found that our hypotheses were partially corroborated: subjects in the *control* and *primed* conditions were significantly less likely to cheat than those in the *loading* and *security* conditions ($\chi^2_3 = 10.676$, $p < 0.014$). This indicates that subjects were more likely to cheat when they had to wait, except when they were told exactly why they had to wait; the label on the progress bar made no observable difference, except when subjects were informed of the danger of PDF viruses and that they were being protected by our software.

Condition	N	0 Pages	1-9 Pages	10-11 Pages
Control	119	3 (2.5%)	12 (10.1%)	104 (87.4%)
Loading	91	4 (4.4%)	20 (22.0%)	67 (73.6%)
Security	87	8 (9.2%)	17 (19.5%)	62 (71.3%)
Primed	103	6 (5.8%)	12 (11.7%)	85 (82.5%)

Table 2: The number of subjects in each condition who read none of the document, some of the document, or all of the document before submitting a response. We considered subjects in the first two categories to have cheated. Recall that we considered those who reached either page ten or eleven as completing the document, since the tenth page was half empty, while the eleventh page explicitly stated it was the end of the document.

We further hypothesized that we would observe two types of cheating: subjects who submit answers before viewing any of the document and subjects who submit answers before reaching the last page (but after opening the document). The former type of cheating happens before subjects experience any types of delays, and therefore should be equally distributed across all of the conditions. Indeed that was the case: a chi-square test indicated no significant differences between the groups with regard to subjects who submitted an answer without ever viewing the document. We therefore decided to remove these subjects from the rest of our analysis, since they did not provide us with any data relevant to our hypotheses. Our results are robust to the point of yielding significance even with these subjects.

We examined the second type of cheating, subjects who submitted answers after only partially reading the document, and found significant differences between the conditions ($\chi^2_3 = 8.619$, $p < 0.035$). Furthermore, we believe that this effect was diminished by the ability to “return” a Mechanical Turk task without receiving credit. Subjects who did not wish to complete the task—but who also did not wish to cheat by entering an arbitrary answer without reading the entire document—had the ability to return the task. While there was not a significant difference between the number of subjects assigned to each condition ($\chi^2_3 = 6.200$, $p < 0.102$), the data in Table 2 shows a disproportionate number of subjects assigned to the *control* and *primed* conditions. Unfortunately, we did not record data from the 55 subjects who returned the task, and therefore we cannot verify whether more subjects were assigned to these conditions and returned them, or if these disproportionate frequencies are due to the random number generator.

Condition	N	Speed Rating	Speed Concerns	Security Concerns
Control	63	4.38	6 (9.5%)	4 (6.3%)
Loading	40	3.53	13 (32.5%)	5 (12.5%)
Security	38	3.71	10 (26.3%)	5 (13.2%)
Primed	46	3.65	12 (26.1%)	14 (30.4%)

Table 3: The differences in survey responses with regards to perceptions of speed and security during the experiment. The columns show the number of respondents in each experimental condition, their average response using a 5-point Likert scale to rate the speed of SuperViewer, and the number of respondents in each condition who explicitly mentioned performance or security concerns.

4.2 Task Performance

We examined the accuracy of subjects’ responses, as well as the amount of time they spent performing the task to determine if there were any differences in their performance based on the conditions. We first examined the accuracy of subjects’ answers, and found no significant differences between the four conditions. However, we did observe that in every condition, cheaters reported significantly different results from those who read the entire document (*control*: $t_{114} = 9.641$, $p < 0.0005$; *loading*: $t_{85} = 10.989$, $p < 0.0005$; *security*: $t_{77} = 10.029$, $p < 0.0005$; *primed*: $t_{95} = 11.516$, $p < 0.0005$). Likewise, we observed that the distance of subjects’ answers from the correct answer was inversely correlated with the number of unique pages they visited ($r = -0.775$, $p < 0.0005$). That is, the more pages the subjects read, the closer their answers were to the correct answer.

We found no significant differences between the experimental conditions when we examined how long subjects spent reading each page. Across all conditions, those who cheated spent significantly less time on the task: $t_{377} = 6.089$, $p < 0.0005$. However, we did notice significant differences based on experimental condition when we examined how many pages subjects viewed. Table 1 depicts the average number of unique pages each participant viewed, as well as the average number of total pages (i.e., counting the same page multiple times if a participant revisited that page).

When we examined the total number of pages subjects visited on average, we found that those in the *control* condition revisited significantly more pages than those in the other conditions ($F_{3,375} = 3.580$, $p < 0.014$). Upon performing post-hoc analysis using Tukey’s adjustment for multiple testing, we found that those in the *control* condition viewed significantly more pages than those in the *loading* ($p < 0.043$) and *security* ($p < 0.029$) conditions. We were particularly surprised to not find a significant contrast between the *control* and *primed* conditions when it came to total pages viewed. Subjects were tolerant of a security delay when they were given a plausible explanation for it, to the point that they were more likely to review their work than subjects in the other two experimental conditions.

4.3 Exit Survey

After subjects completed the experimental portion of this study, we sent them an email offering them an additional \$0.50 in exchange for completing a survey on their opinions of SuperViewer. We received a total of 193 valid responses.

After filtering out subjects who cheated by submitting an answer without ever opening the document, we were left with 187 responses. These respondents claimed to be 125 men and 62 women, with 81% of our respondents holding a college degree or higher. It should be noted that all demographic data was self-reported and unverified, and therefore it may not be representative of reality. Likewise, survey respondents were self-selected from our population of experimental subjects, and therefore may not be representative of the entire population. However, the proportion of cheaters who responded to our survey did not differ significantly from the proportion of cheaters who participated in our experiment. Of these 187 respondents, 27 of them (14.4%) cheated during the experiment by submitting a response after only partially reading the document.

On the first page of our survey, we asked respondents to report their overall impressions of SuperViewer using a 5-point Likert scale. We also asked respondents to rate several factors that contributed to this impression: ease of use, color, look and feel, speed, and security features. When we compared the four different conditions using an ANOVA, we observed a significant difference when it came to perceptions of speed ($F_{3,183} = 8.110$, $p < 0.001$). Upon performing post-hoc analysis using Tukey’s adjustment for multiple testing, we found that people in the *control* condition rated speed significantly higher than those in the *loading* ($p < 0.003$), *security* ($p < 0.036$), and *primed* ($p < 0.011$) conditions. These findings were expected, since those in the *control* condition were not subjected to additional waiting times. The average ratings are displayed in Table 3. Despite the differences in speed, we noticed no difference between the conditions regarding the impact of security features on respondents’ overall opinions of SuperViewer. Over 50% of respondents said that “ease of use” was the primary factor that influenced their overall opinion of SuperViewer, which was consistent across the four conditions.

When we asked subjects what they disliked most about SuperViewer, 40 of them (21.4%) said something about performance or the time it took to load each page:

- *It was very slow.*
- *Pages a little slow to load.*
- *Loading a page with “security features” took an obscene amount of time.*
- *The loading time of the security features.*

We performed a chi-square test to examine whether a dislike for the speed was predominant in any of the experimental conditions. Significantly fewer people in the *control* condition reported speed as being the source of their dislike ($\chi^2_3 = 9.167$, $p < 0.027$). Thus, it is likely that this effect was in response to the presence of the progress bars in all three experimental conditions. The breakdown of these responses are displayed in Table 3.

We asked subjects whether they felt there was a danger viewing the document with SuperViewer and to rate that danger using a 5-point Likert scale. We noticed that 28 of the respondents (15%) said that they had no idea and therefore could not rate the danger. We therefore removed these respondents when we analyzed this question. We observed significant differences between the experimental conditions with regard to subjects’ perceived dangers

($F_{3,155} = 8.636$, $p < 0.034$). Upon performing post-hoc analysis using Tukey’s adjustment for multiple testing, we found that this difference was due to respondents in the *primed* condition rating the danger significantly higher than respondents in the *control* condition ($p < 0.023$). As a follow-up question, we asked respondents to describe any concerns that they had during the experiment. A total of nine subjects in the *control* and *loading* conditions mentioned security concerns, despite the fact that they were not primed to security during the experiment itself, while a total of nineteen subjects in the *security* and *primed* conditions also mentioned specific security concerns ($\chi^2_3 = 12.608$, $p < 0.006$). Using Fisher’s exact test along with the Bonferroni correction to account for multiple testing, we discovered that significantly more people in the *primed* condition mentioned security concerns than those in the *control* condition ($p < 0.0013$). Some of these concerns included:

- *Security is my major concern here. Is it really safe to view PDF?*
- *Wasn’t sure if it contained a virus.*
- *I was a bit concerned when it ran a (sort of) virus scan before displaying the text.*
- *What if there is a virus attached to the PDF file?*

5. DISCUSSION

Our experimental findings were slightly different than what we expected. In the physical world, people must undergo various security mitigations of questionable effectiveness, all the while remaining fairly complaisant. Schneier has written at length about “security theater,” security measures that have no security value other than demonstrating to the public that *something* is being done [28]. Photo identification is checked at office buildings to compare visitors to non-existent watch lists, liquids are banned from airplanes despite evidence that attacks using liquid explosives are impractical, and soldiers with unloaded weapons are placed in prominent public places. Yet the public in general is fairly tolerant of these ineffective security measures, even though they are inconvenient both in terms of time and cost. While the TSA arguably causes more visible inconvenience and delay than most U.S. government agencies, elected representatives do not receive enough constituent complaints for them to actually change policy. In fact, it is likely that investment in technologies such as full body scanners is done mainly to ease perceptions of security, rather than to increase actual security [3]. This lead us to hypothesize that humans may be tolerant of these inconveniences at the mere mention of “security,” whether rational or not.

In this section we explore how our hypotheses compared with the data we collected. We discuss several possible explanations for participants’ behaviors and explain the greater applications for research in this area. Finally, we discuss some of this study’s shortcomings and outline future work in this area.

5.1 Explanations

Our motivation for the contrast between the *security* and *primed* conditions was to examine the role of bounded rationality when people tolerate the security delays. Those in the *security* condition had no rational reason to tolerate the

delay, since no explanation was given other than the ambiguous “security scan” label on the progress bar. Whereas those in the *primed* condition were given a plausible explanation for the security scan. If no significant differences were detected, we hypothesized that bounded rationality would be the case: participants would not need to understand the security explanation to comply. Based on the parallels between our experiment and previous work on soft paternalism and bounded rationality, we did not expect to observe a statistically significant difference in behaviors between these two conditions. We were surprised that this was not the case.

Taken at face value, our results indicate that when given a valid explanation for a security delay, people will tolerate it. While at the same time, without a plausible explanation, people will not tolerate the same delay regardless of whether they are told it is for “security purposes.” We believe there are several possible explanations for why these results diverge from observed behaviors in the physical world.

5.1.1 Habituation

Computer security concerns have grown to prominence in recent years, to the point that even casual users must interact with security mitigations. Users are told to keep antivirus software up to date, to visit only secure websites denoted by a lock icon, and to think critically about what software they install. Yet from users’ perspectives, they see a large quantity of computer security mitigations, but a fairly low attack rate. This calls into question whether these mitigations are worth the cost to users [14]. In the case of SSL warnings, one of the most noticeable user-facing computer security mitigations, the false positives far outnumber the actual positives. Users become habituated to ignoring security warnings because the warnings either do not explain the risks or threat model, or they use jargon such that users do not comprehend them [11, 29]. Thus, users become habituated to many computer security mitigations because they see them so frequently and rarely associate them with consequences.

Due to habituation, users ignore security mitigations when they do not understand the risk. This may explain the lack of differences in behavior when comparing the *loading* and *security* conditions. When users were forced to wait for an arbitrary security check that they did not understand, they did not believe they were in any danger. In the exit survey, 30% of the participants in the *primed* condition mentioned security concerns, almost three times as many as in the *security* condition. Whereas those who mentioned security concerns in the *loading* and *security* conditions were roughly equal. Thus, it is possible that without highlighting a specific threat, users are habituated to computer security messaging. This differs from risk perceptions in the physical world because people can better conceptualize the types of threats without having to have them explicitly stated (i.e., hijacking, bombing, etc.).

5.1.2 Quid Pro Quo

Another possible explanation for the lack of cheating among participants in the *primed* condition is that they felt more obliged to complete the task than participants in the three other conditions. Since these participants were told about the threat model and how our software was protecting them, they may have felt like they owed us something in return,

since we were performing a service on their behalf. We observed that participants in the three experimental conditions all ranked the speed of the program as significantly worse than those in the *control* condition. Obviously this corresponds with seeing the progress bar before viewing each page. However, participants in the *primed* condition, while just as annoyed as those in the other two conditions, were more likely to read through to the end of the document. In fact, these participants were twice as likely to revisit pages in order to review their work. Thus, they expended more effort than those in the two other experimental conditions.

5.1.3 Fear

We believe it is possible that participants who were security primed may have inferred that we are security researchers and therefore thought we may be more likely to detect them cheating. Thus, participants in the *loading* and *security* conditions were more likely to cheat because they were less likely to believe they would be caught. While this is a valid explanation, we believe it is improbable. Previous work conducted on security warnings has found that users are more compliant when they understand the risks of ignoring the warnings [11]. Participants in the *primed* condition were significantly more likely to believe they were being protected from a plausible risk; it is more likely that this sense of risk is what motivated them to comply with the security messaging by continuing to wait. This explanation also does not explain why participants in the *primed* condition were likely to review their work; if they did not cheat simply for fear of being caught, it is likely that they would have done the bare minimum necessary to complete the task.

5.1.4 Sunk Costs

In Section 4, when we observed no significant differences with regard to the amount of time taken, we adjusted participants' completion times to account for the amount of time they had to wait in each condition. That is, we were measuring the amount of time participants spent reading the documents, which did not include the amount of time they spent waiting. We did this because the waiting time was artificially created by us. However, when we factor this waiting time back into each condition, we found significant differences ($F_{3,375} = 3.751, p < 0.011$). Those in the *security* and *primed* conditions spent significantly more time in total than those in the *control* condition.

This does not explain why those in the *security* condition were significantly more likely to cheat than those in the *primed* condition. However, it may explain why those in the *primed* condition did not cheat any more than those in the *control*: participants in the former condition had invested almost 40% more time to complete the task! Thus, the sunk time cost may have dissuaded them from abandoning the task.

5.2 Applications

The study presented in this paper and its results have several immediate applications, both for computer security and public policy.

5.2.1 Security Messaging in Software Systems

Previous work [29] has shown the importance of security warnings in software systems in leading people to adopt secure behavior. The present study shows that, in addition to

making systems more secure, good advance warning systems that clearly explain the rationale for a design choice, also render the system considerably more psychologically acceptable, and make people more likely to tolerate the security choices made for them “under the hood” when these come at a cost.

5.2.2 Alpha Testing Security Features

Security features are best tested by a large number of users, as the multiplication of different use cases across a varied user pool can uncover a number of vulnerabilities that would be harder to observe with a limited amount of testing. The novel contribution of our study is to show that alpha testing (similar to pilot studies in a usability context) is also of utmost importance to gauge the psychological acceptability of security mitigation mechanisms. Performing these types of studies is also extremely cost effective: between our pilots, the experiment, the bonus survey payments, and the fees to Amazon, this study cost under \$200 (beyond the researchers' time). In this manner, scientific results can be used to guide engineering decisions.

As a case in point, consider the relatively complex set of warnings that have to be bypassed to accept a self-signed certificate in Firefox 3. The original beta versions of Firefox 3 contained a 11-step bypass mechanism, which, while marginally increasing security, also utterly annoyed users. Eventually, the bypass mechanism was reduced to a more manageable 4-step process, which is still perceived as too lengthy and impractical given the number of self-signed certificates in circulation [29]. Alpha testing could have helped to spot the problem before Firefox 3 was beta-released, which would have avoided public embarrassment, while leading the Firefox developers to focus on designing a better SSL warning system.

5.2.3 Scareware Defenses

An unfortunate consequence of the results we have obtained is that scareware – the tactic of coercing victims into installing fake anti-viruses, or anti-spyware mechanisms that in fact contain attack code – appears like a very viable strategy for malicious entities. While this result is not surprising, our current work helps quantify the strength of the psychological bias we have to combat when devising defenses against scareware. The more convincing the messaging chosen by the attacker, the more likely the user is to tolerate “odd” behavior from the software installed, so long as the user believes that the behavior is the cost of being protected from a potential risk. Figuring out how to counter this bias opens a whole new avenue of research.

5.2.4 Public Policy Applications

Beyond the software realm, the study seems to confirm that, by calling on people's fears, one may improve the psychological acceptability of any action typically considered as annoying. This shows why Schneier's “security theater” [28] may actually indirectly contribute to security, albeit in a different realm than what it is supposed to originally protect. While the security measures deployed in airports may be very ineffective against people managing to smuggle dangerous materials onboard an airplane, they help pacify the vast majority of the population that has to stand in line, often times in uncomfortable positions. In other words, these security measures, however questionable they may be in terms of

actual protection provided, are likely effective at performing crowd control, which in turn improves the overall security of the environment. More generally, our results corroborate Langer et al.’s results showing that people tend to be compliant with requests when given a rationale, rational or not [18].

5.3 Caveats

During this experiment we found that when participants were provided with a detailed security explanation for the delay, they were significantly more tolerant of the delay than participants who did not see a detailed explanation. This in and of itself does not prove that participants were more tolerant because we displayed a security explanation. It merely shows that without a plausible explanation, participants were less likely to tolerate the delay. It is possible that providing a different detailed explanation unrelated to security may have yielded a similar effect. However, this does indicate that when software developers include security mitigations with associated time costs, users will be more tolerant if they are given a plausible explanation for these time costs.

We collected data from a total of 400 participants in our experiment. Of these, a total of 82 cheated by submitting a response without reaching the end of our document. Twenty-one of these cheaters submitted responses without ever reaching the first page of the document. Upon looking at our web server logs, we discovered that 55 additional participants had begun our task but chose not to submit a response, meaning that a total of 455 people accepted our task. Unfortunately, we cannot determine the conditions to which these 55 were assigned. We do know that all 455 were randomly assigned between the four different conditions. However, it is possible that participants in one of these conditions was more likely to return the task incomplete than participants in other conditions. But given that there were no significant differences regarding how the remaining 400 participants were split between the four groups, we find this explanation unlikely. At the same time, given the lopsided demographic data that we gathered from the exit survey, it is likely that some self-selection bias impacted our study.

We therefore can only speculate about the reasons for accepting the task and not submitting a response. Of the 55, ten participants never even launched SuperViewer. An additional six requested the SuperViewer Flash object, but never viewed the first page of the document. This may be due to Flash incompatibilities or it could be that participants were annoyed by the initial progress bar and closed it before the first page could be displayed. Of the 45 who loaded SuperViewer, the median number of unique pages viewed was four. Six of these participants viewed every page of the document, which begs the question of why they did not submit a response (maybe they forgot?). Without collecting additional data, we cannot know how much the delay or security concerns impacted the decision of these participants to abandon the task.

Finally, there may be concerns about the full effect of the priming on subjects in the *primed* condition. It can be argued that we divided participants into two groups at the beginning of the experiment: security-primed and not security-primed. The other conditions could then be interpreted as subgroups of the latter group. Thus, the effect of

the priming information cannot be separated from the effect of the progress bar. This will need to be addressed in future work to determine the effect of the security priming by itself.

5.4 Future Work

In this paper we highlighted some interesting initial findings with regard to users’ tolerance of security delays. We believe that studies in this area are a crucial missing step in the software development process as well as in the computer security community as whole. A security mitigation may solve a specific security problem, but if users are unwilling to accept the time cost associated with it, it has not solved the problem. However, our study was not without its caveats. We have several future experiments planned to refine our results and to pursue new questions in this area.

5.4.1 Controlling Security Priming

As an immediate follow-up experiment, we expect to repeat this experiment using two new conditions. In the first condition, participants will see the tutorial page from the *primed* condition. However, all the progress bars will be labelled as something unrelated to security (e.g., “rendering fonts,” “performing spellcheck,” etc.). In this manner, we will be able to control for the effect of the security tutorial at the beginning of the experiment. In the second new condition, the tutorial will be changed to highlight the new features that are supposedly causing the security delay. Thus, we will be able to determine whether participants were less likely to cheat because of a plausible security explanation, or if any plausible explanation will suffice.

5.4.2 Determining Maximum Tolerance

In this experiment we showed that those in the *primed* condition were no more likely to cheat than those in the *control* condition. However, differences may still exist with regard to how much of a delay participants will tolerate. We have several experiments planned to examine these upper bounds. In one experiment, we expect to randomly vary the amount of time it takes for the pages to load (though remaining constant on a per-participant basis). This will allow us to use a regression to calculate the upper limit of how long participants are willing to wait in each condition.

In another experiment, we expect to increase the number of pages in the document by several orders of magnitude. By making the document unbearably long, everyone will be forced to abandon the task at some point. We will measure the point at which participants abandon the task in reference to their randomly assigned condition.

5.4.3 Latency vs. Bandwidth

The security mitigations that we modeled in this experiment are high latency in nature: participants were interrupted and had to pause their current task until the mitigation had completed running. Once complete, participants were free to resume their task at full speed until they were interrupted again. Other types of security mitigations are high bandwidth in nature: users are not interrupted, but the speed at which they can perform their tasks is noticeably decreased. We expect to perform experiments in this area as well, in order to further explore the types of slowdowns that are likely to be tolerated.

One particular example is the Tor project. Tor is an anonymous proxy network that uses onion routing [9]. One disadvantage of onion routing is that the level of privacy is directly proportional to the number of hops that packets must traverse. Obviously, this means that privacy comes at a time cost. This creates profound design decisions for Tor’s designers when it comes to specifying default configurations. However, these decisions can be made easier with a better understanding of how much delay users are willing to tolerate in the name of increased privacy.

6. REFERENCES

- [1] A. Acquisti. Privacy in Electronic Commerce and The Economics of Immediate Gratification. In *Proceedings of the 5th ACM Conference on Electronic Commerce*, pages 21–29. ACM, 2004.
- [2] A. Acquisti. Nudging Privacy: The Behavioral Economics of Personal Information. *IEEE Security and Privacy*, 7(6):82–85, 2009.
- [3] E. Berman and L. Heger. Scanners Help Economy by Warding Off Fear of Flying. <http://www.cnn.com/2010/OPINION/02/08/Berman.terrorism.scanners/index.html>, February 8 2010.
- [4] J. Brustoloni and R. Villamarín-Salomón. Improving Security Decisions with Polymorphic and Audited Dialogs. In *Proceedings of the 3rd Symposium on Usable Privacy and Security*, pages 76–85. ACM, 2007.
- [5] J. Condit, M. Harren, S. McPeak, G. C. Necula, and W. Weimer. CCured in The Real World. In *PLDI ’03: Proceedings of the ACM SIGPLAN 2003 Conference on Programming Language Design and Implementation*, pages 232–244, New York, NY, USA, 2003. ACM.
- [6] C. Cowan, C. Pu, D. Maier, J. Walpole, P. Bakke, S. Beattie, A. Grier, P. Wagle, Q. Zhang, and H. Hinton. StackGuard: Automatic adaptive detection and prevention of buffer-overflow attacks. In *Proc. 7th USENIX Security Conference*, pages 63–78, San Antonio, Texas, January 1998.
- [7] L. Cranor and S. Garfinkel. *Security and Usability: Designing secure systems that people can use*. O’Reilly Media, Inc., 2005.
- [8] M. Dalton, H. Kannan, and C. Kozyrakis. Raksha: A flexible information flow architecture for software security. In *ISCA ’07: Proceedings of The 34th Annual International Symposium on Computer Architecture*, pages 482–493, New York, NY, USA, 2007. ACM.
- [9] R. Dingledine, N. Mathewson, and P. Syverson. Tor: The Second-Generation Onion Router. In *Proceedings of the 13th USENIX Security Symposium*, August 2004.
- [10] M. Dowd and A. Sotirov. How to impress girls with browser memory protection bypasses. In *BlackHat Briefings Las Vegas*, 2008.
- [11] S. Egelman, L. F. Cranor, and J. Hong. You’ve been warned: An empirical study of the effectiveness of web browser phishing warnings. In *CHI ’08: Proceeding of The 26th SIGCHI Conference on Human Factors in Computing Systems*, pages 1065–1074, New York, NY, USA, 2008. ACM.
- [12] H. Etoh. GCC Extension for protecting applications from stack-smashing attacks (ProPolice). <http://www.tr1.ibm.com/projects/security/ssp/>, 2003.
- [13] J. Grossklags, N. Christin, and J. Chuang. Secure or insure? A game-theoretic analysis of information security games. In *Proceedings of the 2008 World Wide Web Conference (WWW’08)*, pages 209–218, Beijing, China, Apr. 2008.
- [14] C. Herley. So Long, and No Thanks for The Externalities: The rational rejection of security advice by users. In *NSPW ’09: Proceedings of The 2009 New Security Paradigms Workshop*, pages 133–144, New York, NY, USA, 2009. ACM.
- [15] M. Jakobsson. Experimenting on Mechanical Turk: 5 How Tos. <http://blogs.parc.com/blog/2009/07/experimenting-on-mechanical-turk-5-how-tos/>, July 2009.
- [16] R. W. M. Jones and P. H. J. Kelly. Backwards-Compatible Bounds Checking for Arrays and Pointers in C Programs. In *Distributed Enterprise Applications*, pages 255–283, 1997.
- [17] A. Kittur, E. H. Chi, and B. Suh. Crowdsourcing User Studies with Mechanical Turk. In *CHI ’08: Proceeding of The Twenty-Sixth Annual SIGCHI Conference on Human Factors in Computing Systems*, pages 453–456, New York, NY, USA, 2008. ACM.
- [18] E. Langer, A. Blank, and B. Chanowitz. The Mindlessness of Ostensibly Thoughtful Action: The Role of “Placebic” Information in Interpersonal Interaction. *Journal of Personality and Social Psychology*, 36(6):635–642, 1978.
- [19] G. Loewenstein and E. Haisley. The Economist as Therapist: Methodological ramifications of ‘light’ paternalism. In A. Caplin and A. Schotter, editors, *The Foundations of Positive and Normative Economics: A Handbook*, pages 210–245. Oxford University Press, USA, 2008.
- [20] S. McCamant and G. Morrisett. Evaluating SFI for a CISC architecture. In *15th USENIX Security Symposium*, pages 209–224, Vancouver, BC, Canada, August 2–4, 2006.
- [21] Microsoft Corporation. IE8 Security Part 1: DEP/NX Memory Protection. http://blogs.msdn.com/ie/archive/2008/04/08/ie8-security-part-I_3A00_-dep-nx-memory-protection.aspx, 2008.
- [22] S. Milgram. Behavioral Study of Obedience. *Journal of Abnormal and Social Psychology*, 67:371–378, 1963.
- [23] H. Ozdoganoglu, T. Vijaykumar, C. E. Brodley, B. A. Kuperman, and A. Jalote. SmashGuard: A Hardware Solution to Prevent Security Attacks on the Function Return Address. *IEEE Transactions on Computers*, 55:1271–1285, 2006.
- [24] PaX Project. Address space layout randomization. <http://pageexec.virtualave.net/docs/aslr.txt>, Mar 2003.
- [25] P. Ratanaworabhan, B. Livshits, and B. Zorn. Nozzle: A defense against heap-spraying code injection attacks. In *Proceedings of the Usenix Security Symposium*, August 2009.
- [26] J. Ross, A. Zaldivar, L. Irani, and B. Tomlinson. Who are the Turkers? Worker Demographics in Amazon Mechanical Turk. Technical Report SocialCode-2009-01, University of California, Irvine,

2009.

- [27] B. Schneier. The Psychology of Security. *Communications of the ACM*, 50(5):128, 2007.
- [28] B. Schneier. Is Aviation Security Mostly for Show? <http://www.cnn.com/2009/OPINION/12/29/schneier.air.travel.security.theater/index.html>, December 2009.
- [29] J. Sunshine, S. Egelman, H. Almuhiemedi, N. Atri, and L. F. Cranor. Crying Wolf: An Empirical Study of SSL Warning Effectiveness. In *Proceedings of the 18th USENIX Security Symposium*, 2009.
- [30] R. Thaler and C. Sunstein. *Nudge: Improving decisions about health, wealth, and happiness*. Yale University Press, New Haven and London, 2008.
- [31] J. Tsai, S. Egelman, L. Cranor, and A. Acquisti. The Effect of Online Privacy Information on Purchasing Behavior: An experimental study. *Information Systems Research*, 2010, Forthcoming.
- [32] R. Wahbe, S. Lucco, T. E. Anderson, and S. L. Graham. Efficient software-based fault isolation. In *SOSP '93: Proceedings of The Fourteenth ACM Symposium on Operating Systems Principles*, pages 203–216, New York, NY, USA, 1993. ACM.