# "I did it because I trusted you":
# Challenges with the Study Environment Biasing Participant Behaviours

Andreas Sotirakopoulos, Kirstie Hawkey, Konstantin Beznosov

Laboratory for Education and Research in Secure Systems Engineering
Department of Electrical and Computer Engineering, University of British Columbia
Vancouver, Canada
{andreass,hawkey,beznosov}@ece.ubc.ca

## ABSTRACT

We recently replicated and extended a 2009 study that investigated the effectiveness of SSL warnings. Our experimental design aimed to mitigate some of the limitations of that prior study, including allowing participants to use their web browser of choice and recruiting a more representative user sample. However, during this study we observed and measured a strong bias in participants' behaviour due to the laboratory environment. In this paper we discuss the challenges of observing natural behaviour in a study environment, as well as the challenges of replicating previous studies, given the rapid changes in web technology. Finally, we propose alternatives to traditional laboratory study methodologies that can be considered by the usable security research community when investigating research questions involving sensitive data where trust may influence behaviour.

## Categories and Subject Descriptors

H.4 [**Information Systems Applications**]: Miscellaneous; D.2.8 [**Software Engineering**]: Metrics—*complexity measures, performance measures*

## General Terms

Human Factors, Security, Experimentation

## Keywords

Usable Security, Experimental Design, Study Environment Bias, Study Replication

## 1. INTRODUCTION

Laboratory experiments are often employed in usable security research as they are considered both precise and the best way to control related variables. They have been used

in several studies in order to investigate participant behavior in security contexts without relying on self-reporting of behaviors [2, 8]. The laboratory environment enables experimenters to have better control over the conditions that play a role in the investigated phenomena. It also allows them to observe the participant's behavior first hand, thus avoiding the uncertainty that self-reported data from participants would introduce into the results of the study.

However, there are a number of limitations inherent to laboratory studies. Namely, experimenters have to sacrifice realism and the range of the population to whom the findings can be generalized [6]. These limitations can affect the results or influence the investigated behavior of participants significantly, which may in turn lead to inaccurate conclusions. Usable security researchers [11, 10] have reported that participants in their studies might have altered their behavior in an effort to be "good participants" and to meet the expectations of the experimenter. A number of methods have been considered to mitigate the problems that arise in such an experimental design, including obfuscating the real purpose of the study [11, 10], role-playing [8], and using fake "personal data" to avoid privacy concerns of participants using their own data [11]. However, in many cases these measures have been found to be inadequate [8, 11].

The aim of this paper is to discuss a series of limitations and undesirable effects that were observed during our recent laboratory study on the effectiveness of SSL warnings [9]. Our study was designed to replicate and validate a study conducted at CMU [10], which we will refer to as the CMU study. As in the CMU study, we required participants to perform a series of four tasks; and we observed their reactions to the SSL warnings that were presented to them. After they completed the tasks, we asked them to complete an online questionnaire where we probed their reasoning behind their actions in the study's tasks. As we discuss in the design challenges (Section 2.1), we expanded the CMU design to mitigate some of its limitations.

Although, our preliminary analysis indicates some interesting differences with the results of the CMU study [9], what we find most interesting is the reasoning participants gave for their decision to ignore the warnings. A substantial number of them claimed that they ignored the warnings because of the safety they felt due to the study environment. Furthermore, due to the design of our recruitment phase, we were able to observe that a large portion of the sample population was reluctant to participate in the study as it re-

quired the use of private data (i.e., their bank credentials). We believe our results have can inform the usable security research community about the limited ability to draw conclusions from participants' actions during laboratory studies. We found that the advantages of the laboratory environment were countered by its limitations. As similar experimental designs are common in usable security studies, we argue that in the case of behavioural research in the context of computer security, alternatives to a lab study methodology or the use of complementary methodologies should be considered. We next present our recruiting and experimental procedures in detail. We then discuss the observations that we have made about the representativeness of our participant population, the impact of the study environment on our findings, and the differences between self-reported data and observed actions.

## 2. EXPERIMENTAL DESIGN

Our experiment was designed so that it would mitigate some of the limitations of the CMU experimental design, both these acknowledged in the CMU study and those which we felt should be addressed. The first limitation we identified was that participants in the CMU study were drawn almost exclusively from the CMU student body. A second limitation was that their participants were randomly assigned to the browsers investigated, which might have caused them to alter their normal behaviour and become more cautious about SSL warnings as the warning interface was unfamiliar to them. Thirdly, in the CMU study, the custom warnings designed for Internet Explorer 7 (IE7) were radically different in colors, wording and layout from the native IE7 warnings (Figure 1). We believe that this also might have contributed to participants being surprised and eliciting a more cautious reaction to the warnings. In an effort to mitigate these limitations, we recruited participants from the broader Vancouver population instead of limiting ourselves to UBC students. We assigned users to our conditions according to the browser they normally used; and we redesigned the custom SSL warnings that were presented to the users, keeping the layout similar to the native warning in an effort to limit the surprise effect a previously unseen browser interface and warning would have (Figure 2).

Our study was a between subjects experiment with four conditions based on the warning presented and browser used by the participant. The four conditions were the following:

- Firefox 3.5 browser presenting its native SSL Warning

- Firefox 3.5 browser presenting an SSL warning designed by us

- Internet Explorer 7 presenting its native SSL Warning

- Internet Explorer 7 presenting an SSL warning designed by us

We did not replicate the CMU study's Internet Explorer 6 (IE6) condition, because we could not recruit sufficient number of participants that used IE 6 in their everyday life. Finally, we did not replicate the custom multi-page warning condition present in the CMU study. Instead we substituted it with a Firefox 3.5 (FF3) condition, in which the custom warning also retained the same layout and changed the colors and wording. We did this because we wanted
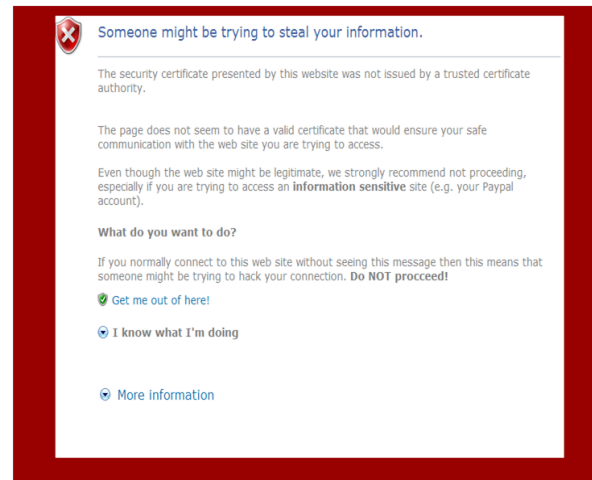


**Figure 1: CMU custom warning**



**Figure 2: Our custom warning for IE7**

to investigate if there are differences between security practices between IE and FF users. By using the same colors and wording but maintaining each browser's layout, we would be able to directly compare reactions to our custom warnings between users of these browsers.

We later added a fifth condition to our study as we realized that because we had changed two variables from the CMU study conditions (i.e., breadth of population, use of regular web browser), we would not be able to draw conclusions about the cause of any differences that we might observe when analyzing our data. In the fifth condition, we also used a broad population (i.e., not only university students and staff), but we asked them to use IE7 regardless of their what browser they would normally use. This condition should enable us to attribute whether any differences in the reactions to the warnings are due to the population sample used or the choice of web browser.

We next describe the design challenges we faced, the steps we took to address them. We then describe our recruitment process, the study tasks, and the exit questionnaire and discuss any differences between our study protocol and the CMU study protocol.

## 2.1 Design Challenges

There were several broad challenges that we had to consider as we designed our study. We did not want our participants to be primed for security so we did not want to reveal our study purpose. Also, we wanted to avoid them being surprised by being presented with an unfamiliar interface that might cause a more cautious than normal reaction to the SSL warnings encountered. In addition, we wanted to recruit participants with a diverse demographic background. As described next, we also faced several challenges as a result of this being a replication of a previous study.

### 2.1.1 Challenges in replication a previous study

During our efforts to replicate the CMU study, we came across challenges due to the rapidly changing web environment in which the SSL warnings are raised. In [3] it is argued that replication of past studies in a web environment might face contextual challenges due to the fact that the web constantly changes and evolves. In point of fact, when we began to design our experiment, we realized that we had to deal with a completely different contextual environment than the one in which the CMU study was conducted. In 2008, when the CMU study was designed and conducted, FF3 had only recently been released (June 2008). Consequently, users were much less accustomed to its interface and warnings; and, therefore, might have been surprised when they faced such a new interface. We speculate that this is one reason for the effectiveness of the FF3 warning in the CMU study, as the FF3 warning had changed layout completely and become rather complex in the latest version released prior to the study. We also realized that we would be unable to recruit participants that used FF2 in their every day online tasks, as after almost 2 years, very few users of Firefox had not upgraded to FF3.

The lesson we learned is that, especially in a web environment, replication studies have to face potentially quite different conditions from the original study. These differences might affect the ability to directly compare findings with the earlier study, even when conducted after a relatively short period of time.

## 2.2 Steps to Address Design Challenges

In order not to prime participants for security, we advertised our study as one that investigates challenges people face when retrieving information online. The same was written in the consent form that participants read and signed prior to the experiment.

In addition, for the purpose recruiting a broader population, we not only posted our flyers around campus but also at community centers around Vancouver; and we advertised the study on Craigslist. In an effort to recruit, as representative a sample as possible, the only conditions we imposed were that participants should be familiar with one of the two browsers that are under investigation (IE7, FF3) and that they have an online banking account. Due to the difficulty of recruiting non-student participants to take part in a study on the UBC campus, we eventually had to raise the honorarium offered to participants from $10 to $20 so as to meet our recruitment goals.

## 2.3 Recruitment

We adopted a three step approach for our recruitment. First we advertised the study using flyers around UBC campus, Vancouver community centers and Craigslist. In our advertisements, we mentioned that users will seek information from various sources like Google, online banking sites, and online shopping sites; but we did not reveal that they will be asked to use personal information to do so. When potential participants contacted us via email or telephone, we set a date and time for a session via email. After we had a session arranged, we sent a final email with the consent form attached, as required by the UBC ethics board, and details about the location of the study. It was only then that we revealed that they will have to retrieve information from various on-line sources, including their bank's online system; and for that purpose, they had to have an account with one of the major Canadian banks and should remember their bank credentials. We not only wanted to ensure that participants had an online banking account, but also wanted to ensure that they would have their client card number with them (or at least remember it) as this is used in most cases as their username.

The CMU study's recruitment process advertised the study as a "usability of information sources study." Also only people that were customers of one particular bank were used as participants in the CMU study. A screening process was administered that required participants to have used search engines and performed an online purchase in the last year. In contrast, we advertised our study as one that seeks to investigate the challenges users face when retrieving information online and that we seek to identify "difficulties people are facing when trying to compete every day tasks online (e.g. search on google.com for information, online banking, online shopping)". We omitted the screening survey done in the CMU study because we wanted to impose a minimum of requirements for our participants as we were aiming for a broad population both in terms of occupation and age. We felt that older participants might be relatively unfamiliar with tasks like online shopping. In retrospect, this change might have affected how well we were able to obfuscate the real purpose of our study. However, we have no evidence that this actually happened.

Finally we purposely set a date and time for the study prior to revealing that real account information would be used by participants during the study. We hypothesized that if participants were concerned with the privacy of their information, they would explicitly state that when they would cancel the session. Although we had some participants that did not come to their session and did not provide any reasons for missing it, our hypothesis was confirmed in several cases.

## 2.4 Tasks

When a participant arrived, the experimenter gave him a copy of the consent form and asked him to sign it. The experimenter then gave a detailed overview of the study proceedings without revealing the real purpose of the study. Four tasks were then presented to the participant and he was asked to review them and ask any question that he might have. Each task asked for a piece of information and included a primary source and a secondary source that would enable the participant to retrieve the information. This was a feature of the CMU study, which aimed to mitigate the task focus effect that has been observed in similar studies.

The first task asked participants to retrieve the surface area of Greece using Google.com as a primary source and

Ask.com as the secondary one. The second task asked participants to retrieve the last two digits of their account balance using either the online banking system as the primary source or the phone banking system of their bank as the secondary source. The third task asked participants to locate the price of the hardcover edition of the book Freaknomics using either Amazon.com as the primary source or Barnes and Noble as the secondary source. The fourth task asked participants to create a new email account in order to register with tripadvisor.com, using either Hotmail.com as the primary site or Yahoo.com as the secondary site. The first and third of the tasks were dummy tasks that were there only to obfuscate the real purpose of the study and reinforce to the participants' belief that this study was not about warnings. Half of the participants had the order of the bank and email tasks (i.e., the warning tasks) swapped so that we would control for any order effects in the warnings presented.

In the CMU study, the email task was not used at all. Their fourth task asked participants to use the CMU online library catalog or alternatively the library phone number to retrieve the call number of a book. As we wanted to recruit participants from outside UBC, we could not design a similar task. We opted for the email task as described above, which we feel serves well as a task with a relatively low risk for personal information exposure for the participant.

In our study, the experimenter did not help the participants during the study (although many asked for help while performing the tasks, including the dummy tasks), but did not deny that he was part of the research team. However, in subsequent discussions with of one of the CMU researchers (S. Egelman, personal communication, March 31, 2010), we found out that in the CMU study, the researcher who administered the study pretended that he had no connection with the research other than getting paid to sit in the room with the participants and just read the script to them. This subtlety was not reported in the paper describing the CMU study [10].

## 2.5   Exit questionnaire

After the completion of the four tasks, participants were directed to an online SurveyMonkey questionnaire. Similarly to the CMU study, the questionnaire asked 45 questions in six categories. The first set of questions asked about participants' understanding of and reaction to the bank warning in the study. The second question asked the same questions about the Hotmail warning. The third set asked questions to gauge their general understanding of certificates and invalid certificate warnings. The fourth set gauged participants' prior exposure to identity theft and other cyber threats. The fifth set, asked them about their technical experience, including their experience with computer security. Finally, the sixth set asked general demographic questions like age, gender and education level.

We kept most of the same questions from the CMU study, but we added some in order to further investigate points of interest (e.g., if the participant would perform differently if using his own PC). While participants responded to the questions presented to them, we asked them to elaborate on their answers where appropriate. For example, in the question asking them if they would otherwise react to the warning if it was their own PC, participants had to choose between "yes" and "no". We then asked them why and if

their answer would change if the PC was one in an internet cafe.

After the questionnaire was completed, we debriefed participants and revealed the true purpose of our study. We also explained to them the utility of SSL and the SSL warnings.

## 3.   DISCUSSION

The overall aim of our experiment is to investigate computer security behavior in the context of SSL warnings. While the study is still underway and we are not yet ready to report results, it is evident that there is a significant impact of the laboratory environment, which introduces bias and uncertainty. This impacts not only the results gathered and their quality, but also on the profile of participants that took part in our study. As we will discuss in more detail, this impact is due to systematic limitations of the experimental method; therefore, the key points discussed are applicable to other studies in the field that have a similar experimental design and overall goals.

## 3.1   Representative participant population

The participant pool is of outmost importance for any study that requires a representative population sample. Without a representative sample, it is hard to generalize results or draw conclusions, especially when what is investigated is the behavior or views of individuals. In our study, we have a broad population in terms of age, education, and occupation; however, in reality, our participant population is skewed. Participants are not randomly recruited, rather they volunteer; and it became clear during recruitment that many security concerned people opted not to participate.

Due to our recruiting method, we were able to establish communication with potential participants before making it clear that they would be required to use sensitive, personal information (i.e., log into their actual online banking site). This allowed us to observe the fact that the most security aware or cautious individuals decided, either on their own or because of an advice by someone in their environment (e.g., spouse), not to take part in the study. To the present date, 68 participants have taken part in our study from 125 initial contacts. However, 10% (14/125) of those who contacted us refused to participate after finding out about the banking task. This refusal occurred even though we stressed that no information would be recorded during the study and we sent the consent form, as an attachment to the same email, which may have provided an additional sense of safety and security.

This raises a concern about the statistical validity of the recruited sample when sensitive information must be used by a participant and he has prior knowledge of that. The problem is that a considerable percentage of the potential participants will not take part in the study out of fear of leakage of their information. The systematic error introduced here is due to the fact that the users who take part in studies similar to our own fit a certain behavioral profile and so conclusions can be drawn only for this behavioral profile. This potentially affects the reported severity of the problem under investigation as we are missing recruiting as participants those users who do the "right thing" (i.e., keep their information safe and private). As a result the generalizability of the conclusions drawn on the results gathered is degraded.

| Reason for My Reaction to the Bank Warning | |
|---|---|
| It is a study | 33% |
| Calling the bank is time consuming | 15% |
| I wanted to complete the task | 13% |
| I am used to the warning | 10% |
| I trust my bank's web site | 25% |
| Other | 4% |

Table 1: Participants' responses, in the online survey, when asked why they ignored the warning at their bank's web site

## 3.2 Impact of study environment

Although the authors of the CMU study speculated that the study environment might had affected their results, they do not report any relevant data on this. During the exit questionnaire, we asked participants in the form of an open question why they chose to ignore or heed the SSL warnings with which they were presented. This allowed us to measure through self-report data interesting findings about the reasons that participants had for their action upon seeing the SSL warning at the bank login web site (Table 1). For those who ignored the warning at the bank site, one third explicitly stated that it was the study itself (i.e., being sanctioned by UBC and having approval by the ethics board) that made them trust the procedure and the experimental setup and ignore the warning in order to enter their personal information into the site. Furthermore, most of those that said that it was the environment who played the most significant role in their decision claimed that they would do otherwise if this was not a laboratory experiment and they saw that warning at a public PC or network or even their own computer.

Another 13% of participants claimed that they ignored the warnings because they wanted to complete the task. If considered conservatively, a portion of this percentage can be interpreted as a task focus effect (i.e., the participant is continuing to the web site because he feels that he should do as asked) [7].

These two types of response make us question the very utility of laboratory study designs in usable security. Although this is not, by any means, new knowledge in the field of behavioral studies; it is the first time to the best knowledge of the authors, that concrete data have been collected in a usable security study suggesting inadequacy of the laboratory study as a tool for this purpose. Even if measures are taken to mitigate this issue, in our opinion it is very hard, if not impossible, to make sure that the researcher has been successful.

We argue that using laboratory studies as an experimental methodology contradicts popular computer security advice. In order to have users as participants, we essentially ask them to perform in public and potentially unsafe environments those actions that they have been told/trained to perform in private and safe environments. Consecutively, we have a systematic error introduced in the study design and potential participant behavioral profile. Namely we either have users, as participants, who are by nature prone to unsafe behavior or we have participants that put quite justified good faith in the study environment due to the reassurance of ethics boards and prestigious institutions. In either case, the conclusions drawn from the results collected

under such circumstances can hardly be considered reliable and representative of actual user practices. It is the authors' belief that although laboratory user studies are invaluable in usability research, they are unfitting for the purpose of research in usable security when the investigation includes the need of participants to use personal data in the study environment. The systematic error introduced by the study environment seems to be a fundamental one.

Measures (e.g., adding tasks to conceal the real purpose of the study) are often inadequate to mitigate this error. It might be possible to successfully conceal the real purpose of an experiment by adding more irrelevant tasks and devise elaborate scenarios to create a greater sense of realism. However, adding too many tasks could render the study infeasible due to time constraints and scenarios do not seem to be successful in creating a realistic environment [8]. However, proceeding without taking into consideration the ethical implications and constraints that a study should have could prove equally problematic, as happened in the case of [4].

We suggest that it is conceivable that alternatives to a regular laboratory study may have to be sought in cases where we seek to investigate human behavior in computer security by having participants use their private data. In an effort to maintain a sense of realism without compromising the ethical aspect of a study, Jakobsson et al in [5] designed an experiment that took into consideration both realism and ethical concern. Moreover, researchers could devise experiments in lieu of laboratory studies that use methods to log users' behavior while they unknowingly perform everyday actions. Then, after asking for the users' consent, they could analyze these logs (e.g., contacting the IT department of a big corporation and installing a proxy that generates SSL warnings on particular sites and logging the reaction). Such post-hoc designs have been used by the HCI research community in investigating web use [1]. Although we are aware of the challenges in terms of research ethics and acquiring the cooperation of organizations and individuals, it is our belief that the benefits in terms of the reliability of results will be substantial. We propose that additional attention and effort should be dedicated towards such experimental designs and away from laboratory studies in future usable security research that investigates user behaviors.

## 3.3 Differences in self-reported and observed actions

During our exit questionnaire we presented participants with a screenshot of a Firefox 3.0 SSL warning raised due to a self signed certificate. We asked what they would do if they saw this warning prior to entering a web site (e.g., www.example.com) and provided a multiple choice question with four choices: "I would leave the web site", "I would proceed if the web site was not information sensitive", "I would proceed to the web site", "Other, please explain". As shown in Table 2, the majority of participants claimed that they would leave the web site or would proceed if it was not information sensitive. This is very different from the actual behaviours that we observed when participants were presented with an SSL warning during the banking and Hotmail tasks; in that case, the majority of participants ignored the warnings and proceeded, even when facing custom warnings that had intense colours and strong wording.

This question was asked in the late stages of the exit ques-

| Hypothesized Action for FF3 Warning | |
|---|---|
| I would proceed ignoring the warning | 14% |
| I would proceed if the site was not information sensitive | 28% |
| I would leave the web site | 43% |
| Other, please explain | 15% |

**Table 2: Self reporting of participants' normal action when presented with a screenshot of the Firefox 3.0 self signed certificate warning while trying to reach a hypothetical web site**

tionnaire, so it is quite possible that participants understood the actual purpose of the study (i.e., their reaction to SSL warnings) by the time they answered the question. In that case, they may have been biased to provide the "correct" response rather than accurately report their usual behaviour.

It would be interesting to study further the reasons behind such differences. We are aware of the problems of self-reporting in terms of reliability, but we also have to take into account the bias the laboratory environment brings into the experiment. In order to conclude which method produces more reliable results for the context under investigation, we would like to design an experiment that would ask participants in a survey for qualitative responses on what their reaction would be if presented with a security feature. At another time (either before or after), their actions could be observed in a natural setting as they are presented with warnings during their normal tasks. Although we are aware of the challenges that such an experimental design would involve (e.g., ethics approval, lack of controlled environment), it is conceivable that the experiment could be conducted outside a lab environment with minimal or no awareness on behalf of the participants. This way the laboratory bias would be mitigated, and it would be clear to the researcher the reliability of the survey responses.

## 4.  CONCLUSION

We presented our experimental design used to investigate the effectiveness of SSL warnings and to validate the findings the prior CMU study. We took certain measures to address the limitations of the former experimental design. Although this is a work in progress and we intend on extending it to investigate the problem further, our preliminary analysis has revealed differences between our study and the CMU study under validation. Furthermore, it also raised concerns about the limitations of laboratory studies when used to do usable security research on human behaviors.

It is our belief that the aforementioned limitations of the particular experimental method are applicable to other studies in the field of usable security. The reluctance of security concerned people to take part in the our study that we observed during participant recruitment raises concerns about the ability of such studies to accurately and reliably draw conclusions about security practices and user behavior of the general population. Finally, we proposed alternative study methodologies that might be free of the fundamental errors discussed and could yield more reliable results. Such methodologies might complement current experimental designs in order to mitigate the limitations inherent in any design approach.

## 5.  REFERENCES

[1] A. Cockburn and B. McKenzie. What do web users do? an empirical analysis of web use. *Int. J. Human-Computer Studies*, 54:903–922, 2001.

[2] 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: Proc. of the SIGCHI conf. on Human factors in Computing Systems*, pages 1065–1074, New York, NY, USA, 2008. ACM.

[3] K. Hawkey and M. Kellar. *Handbook of Research on Web Log Analysis*, chapter XI, Recommendations for Reporting Web Usage Studies, pages 179–202. IGI Global, 2009.

[4] T. N. Jagatic, N. A. Johnson, M. Jakobsson, and F. Menczer. Social phishing. *Commun. ACM*, 50(10):94–100, 2007.

[5] M. Jakobsson and J. Ratkiewicz. Designing ethical phishing experiments: a study of (rot13) ronl query features. In *WWW '06: Proceedings of the 15th international conference on World Wide Web*, pages 513–522, New York, NY, USA, 2006. ACM.

[6] J. E. McGrath. Methodology matters: doing research in the behavioral and social sciences. *Human-computer interaction: toward the year 2000*, pages 152–169, 1995. Morgan Kaufmann Publishers Inc.

[7] A. Patrick. Commentary on research on new security indicators - essay. http://www.andrewpatrick.ca/essays/commentary-on-research-on-new-security-indicators, 2007.

[8] S. E. Schechter, R. Dhamija, A. Ozment, and I. Fischer. The emperor's new security indicators. In *Proceedings of the 2007 IEEE Symposium on Security and Privacy*, pages 51–65, Washington, DC, USA, 2007. IEEE Computer Society.

[9] A. Sotirakopoulos, K. Hawkey, and K. Beznosov. Poster: Validating and extending a study on the effectiveness of ssl warnings. Poster at Symposium on Usable Privacy and Security, 2010.

[10] J. Sunshine, S. Egelman, H. Almuhimedi, N. Atri, and L. F. Cranor. Crying Wolf: An empirical study of SSL warning effectiveness. In *Proceedings of 18th USENIX Security Symposium*, pages 399–432, 2009.

[11] T. Whalen and K. M. Inkpen. Gathering evidence: use of visual security cues in web browsers. In *Graphics Interface*, pages 137–144. Canadian Human-Computer Communications Society, 2005.