

# Building Robust Authentication Systems With Activity-Based Personal Questions \*

Anitra Babic  
Computer Science Department  
Chestnut Hill College, Philadelphia  
BabicA@chc.edu

Danfeng Yao  
Computer Science Department  
Rutgers University, New Brunswick  
danfeng@cs.rutgers.edu

Huijun Xiong  
Computer Science Department  
Rutgers University, New Brunswick  
huijun@cs.rutgers.edu

Liviu Iftode  
Computer Science Department  
Rutgers University, New Brunswick  
iftode@cs.rutgers.edu

## ABSTRACT

A recent study found that the widely-used secret questions for Web authentication can easily be guessed. The study focused on making secret questions easier to remember for the user and harder to break by others. Our approach is authentication through the use of an individual's personal and dynamic Internet activities. We hypothesize that frequently-changing secret questions will be hard for attackers to guess. We propose three major categories of questions that are based off of user activities: network activities (e.g., browsing history, emails); physical events (e.g., planned meetings, calendar items); conceptual opinions (e.g., opinions as derived from browsing, emails). Our preliminary results are encouraging and show that this new direction is promising.

To improve the usability, in particular nonintrusiveness, of such a dynamic secret-question system, we also describe a concrete client-server architecture and security model for automating our authentication systems through utilizing existing artificial intelligent techniques.

## Categories and Subject Descriptors

D.4.6 [OPERATING SYSTEMS]: Security and Protection-Authentication

## General Terms

Security

## Keywords

Personal questions, authentication, activity, opinion, security, usability

---

\*This work has been supported in part by DIMACS REU grant, NSF grant CCF-0728937, CNS-0831186, CNS-0831268 and the Rutgers University Computing Coordination Council Pervasive Computing Initiative Grant.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

*SafeConfig'09*, November 9, 2009, Chicago, Illinois, USA.  
Copyright 2009 ACM 978-1-60558-778-3/09/11 ...\$10.00.

## 1. INTRODUCTION

Authentication is an important aspect of a secure system, where a user proves her identity by revealing the fact that she possesses certain secrets or objects, such as passwords, private keys, physical tokens, biometrics, or answers to some personal questions. Different types of authentication mechanisms provide different security and assurance guarantees, and various authentication mechanisms can be combined to improve the robustness of the system. Stronger authentication systems may be harder to use than weaker ones due to the complexity involved in the cryptographic protocols. For example, non-security savvy users may find the PGP 5.0 encryption tool not intuitive to use [8]. Personal questions are a weaker yet more usable form of authentication method, and are widely used on the Internet. At the initiation, a user enters her secret answers to a set of questions, which are stored by the server. When a user is challenged later by the server, her answers need to match the stored ones. Most popular Web mail providers, e.g., AOL and Google, rely on personal questions as the secondary authentication secrets for resetting account passwords.

However, a recent study found that the widely-used secret questions for Web authentication can be guessed [7]. In that study, participants were asked to answer personal questions from four major Web mail providers (AOL, Google, Microsoft, and Yahoo!), and then their acquaintances are asked to guess their answers. Both security and usability were evaluated. The authors found that acquaintances of participants were able to guess 10% of their answers [7]. 13% of answers could be guessed within five attempts by guessing the most popular answers of other participants. Some questions were too hard for participants to remember in a long-term, as participants forgot 20% of their own answers within six months. These negative results indicate that there is an urgent need to investigate new approaches for usable and robust question-based authentication systems that leverage personal knowledge.

In this paper, we focus on making personal questions easier to remember for the user and harder to break by others. Our unique approach is to design activity-based short-lived and dynamic authentication questions by utilizing of a user's personal and dynamic Internet activities. We hypothesize that short-lived and frequent-changing questions will be harder for attackers to guess, and easier for a user to remember. We propose three major categories of questions that are based off of user activities: network activities (e.g., browsing history, emails); physical events (e.g., planned meetings, calendar items); conceptual opinions (e.g., opinions as derived from browsing, and emails). We give more details and ex-

amples on our activity-based personal questions, and our preliminary evaluation results are also presented.

In existing secret question based authentication systems, the correct answers need to be specified by the user prior to the challenge phase. Thus, intuitively, short-lived questions may require frequent manual updates from the users in order to install their correct answers. Such a requirement may make the system feel intrusive to the users. To improve the nonintrusiveness of our proposed authentication mechanism, we design an automatic and secure architecture to automatically extract answers from the users' Web activities without any user participation, e.g., the *correct* answer to question *Who was the last person that you sent mail to today?* can be automatically extracted from Yahoo! mail server *without* the user's manual update. Such a dynamic activity-based authentication system can be deployed as a client-server architecture and utilizes existing artificial intelligent techniques such as opinion extraction methods [1]. We also describe the security models for automating our authentication systems.

**A Use Scenario** Our proposed activity-based personal questions can be used in several situations, one such use being in congruence with an email system. When a user forgets his/her password or wishes to reset password, the email system can ask secondary authentication questions based off of events that occurred within the user's email, no matter which machine the user is on. With user information being held on a server, the user's location does not matter. So long as the user can recall her account name, activity based personal questions can be generated using information from that account, a system that is currently utilized by most email providers such as Yahoo! and Gmail. Like these currently implemented systems, if the user cannot answer the generated activity based personal question they will be unable to access the account or their password reset will not occur, depending on what the user was attempting to do.

With this scenario in mind it is easy to see activity based personal questions being easily transferable to any online account. Provided there are logs of user activity and the user is able to provide their account name all manner of questions encompassing the principles of our activity based questions can be generated that only the user would be able to answer.

We note that our solution does not create new privacy vulnerabilities, because the server leverages its existing transaction logs and stored data about users to generate activity-based authentication questions, without requiring additional user information. As the server is responsible for extracting challenge questions, there is no need to install any client-side software, and thus a user can use any computer without constraints.

The rest of the paper is organized as follows. We first describe the design principles, categories, and examples of our new activity-based authentication questions. Section 3 gives the architecture and security model of an automatic authentication system. We describe our preliminary user study and its results in Section 4. Conclusions and future work are given in Section 5.

## 2. ACTIVITY BASED AUTHENTICATION QUESTIONS

To improve robustness and usability of personal authentication systems, we specify four main requirements that our activity-based questions need to meet. These requirements encompass areas such as:

- **Secrecy:** The correct answers of challenge questions should be hard to guess by attackers.
- **Memorability:** The user only has to recall their most recent

events and network activity in order to answer the generated questions.

- **Non-intrusiveness:** The system should have the potential to run in the background with the computer both updating answers and generating questions.
- **Adaptability:** All the question and corresponding answers should be produced automatically and refreshed periodically.

Secrecy is required in order to ensure that only the legitimate user knows the answer to her own activity based questions. However, within a previous study users were proven to have a low rate of recollection when it came to recalling their own answers to current authentication questions [7]. Thus, the answers to our activity based questions are required to be memorable as well secure because of these results. Our solution is to design short-lived questions based on recent activities of a user. The nonintrusiveness requirement is to eliminate or reduce the need of manually specifying correct answers by a user. We describe ways to achieve nonintrusiveness in Section 3. Last but not the least, adaptivity ensures the freshness of the challenge questions. Our activity based security questions were formed around these requirements, each question falling into one of three categories: network activity, physical events, and conceptual opinions. Some sample-questions are given in Table 1.

**Network Activity** We propose questions that fall into the category of network activity monitor and focus on the user's online activity. With questions that range from *what was the last website you visited* to *how large was the last email you sent out* they focus on the size, type, history, and content of user network activity. This makes questions not only more memorable but also gives users a degree of security. Browsing patterns and habits are usually personal, so that users rarely discuss them in depth with others. User browsing habits vary from person to person and from day to day, a new IP address or URL occasionally being visited or no browsing activity occurring during certain times. However, as we will discuss later, the overall security of these questions will also depend on the individual and on the popularity of the sites that they are visiting.

**Physical Events** The second category of activity based authentication questions focuses on physical user events. This is achieved through information gathered from emails, virtual calendars, social networks, and other planner-like programs. When activities are entered into the calendar or an invitation is accepted via email, a question can be generated from the event as well as its related answer. For example, questions about who the user will be meeting next Monday at 5pm or where the meeting on Thursday will be located are activity-based authentication questions that have an element of secrecy.

However, this category of questions is not without its faults. The secrecy is more relative to how many people are attending the event and how many people are in tune with the user's schedule, giving these questions varying degrees of secrecy. Consequently, while these questions might not be as secure, they are easier for the user to remember due to their being drawn from actual events the user is already consciously trying to remember.

**Conceptual Opinions** Opinions can be extracted from our personal correspondence and through the sites visited while using the Internet. Someone who is of a certain viewpoint will generally visit and read articles with similar sympathies. By analyzing browsing history and email content, it is possible to extract a user's opinion. Breck *et al.* have done work in this area with findings that show the plausibility of having a machine extract text that expresses an opinion and perform an analysis on the text [1]. With this in mind,

questions such as *what is your opinion on North Korean Nuclear testing* can be asked and answered by user.

These questions can be easily adapted to learn the user’s opinion on current events and allow for easy user memorability as well. A personal opinion is arguably hard to forget. However, there is a possibility that opinion-based personal questions may be vulnerable to random guessing attacks. This problem can be mitigated by requiring users to correctly answer  $k$  questions. This method can increase the accuracy of opinion-based questions and make it more difficult for attackers to guess correct answers. For example, if there are three choices (yes, no, and neutral), then the probability of correctly guessing all  $k$  questions is  $(1/3)^k$ , assuming equal likelihood of each choice for all questions.

### 3. SYSTEM DESIGN

In this section, we describe the architectural design of an activity-based authentication system that is capable of automatically generating challenge questions and corresponding correct answers from user’s activity logs. Because of our adaptivity and memorability requirements, the authentication system needs to update and refresh the questions frequently. Further detail about our implementation plan on a proof-of-concept prototype can be found in Section 5.

Our model is a client-server architecture, as in the traditional question-based authentication systems. Figure 1 shows a schematic drawing of our design. The server utilizes the logged user-transaction data to extract personal questions and corresponding answers about an individual. In our security model, we assume that the server is not compromised by malicious software, which means that all application data is secure on the server. In addition, the communication between server and client requires SSL and is assumed to be confidential, and the transaction data during communication is secure.

The authentication service in our system involves two phases, *setup* and *authentication*. During the setup phase,  $\langle$ question, answer $\rangle$  tuples are automatically generated from user’s daily activity data on the server. The tuples are stored in a database as a basis source of secret questions for our automatic authentication. The challenge questions are dynamically generated based on recent activities of a user and are short-lived to ensure both memorability and security. In the authentication phase, e.g., when the user sends a request for retrieving a password, our system presents the user with secret questions generated in *setup* phase, and verifies the correctness of user’s answer by comparing with the automatically extracted answer. The architecture of such a system consists of four components: *Preprocessor*, *Parser*, *Question Generator*, and *Authenticator*. Details of each component are described next.

- *Preprocessor* accesses the user’s activity data logged on the server, and truncates them into lists of small plain text as raw data, which are then passed down to the *Parser*. Different Web services have different formats of user’s activity data, for example, e-commerce server keeps each user’s transaction, while email server stores user’s emails. As a consequence, *Preprocessor* will apply varying trim policy to original activity data in the server, such as transaction-based trim or email-based trim.
- *Parser* interprets the semantic meaning of user’s activity data and converts it into an annotated form. It takes the raw data obtained by *Preprocessor*, extracts activity-related fields, and then stores them into different kinds of tokens with specific types. In the email systems, our system extracts time, sender,

receiver, email title and body, and inserts them into certain type of token according to their semantics. Given an email’s title: *Running in Lincoln Park this Saturday’s afternoon*, this text will be segmented by the *Parser* into different tokens such as  $\langle$ behavior: run $\rangle$ ,  $\langle$ place: Lincoln Park $\rangle$ ,  $\langle$ time: Saturday’s afternoon $\rangle$ .

- *Question Generator* correlates meta data passed down by *Parser* with tags, produces relevant questions and answers, and stores them into a database. The tags in our system include *Who/Whom*, *When*, *Where*, *Who*, *What*, *How many/much*. E.g., token  $\langle$ behavior $\rangle$  is given tag *What*,  $\langle$ time $\rangle$  with *When*,  $\langle$ place $\rangle$  with *Where*. *Question Generator* then produces questions in a natural-language form asking for data related to certain tag(s). Existing tools on natural-language processing (NLP) can be incorporated to realize the question-generation process [6].
- *Authenticator* issues freshly-generated activity-based authentication questions to the user. It then semantically interprets the natural-language based answers from users, and compares them with the correct answers stored in the database. For example, the semantic interpretation requires the *Authenticator* to recognize that the answer of “jogging” is equivalent to “running” in a proper context.

Our model is general and can be deployed on servers that provide network related services for users. The server logs in our model provide information in which questions can be generated. The server may be an email/calendar server or an e-commerce server for on-line shopping and banking. Because the server leverages its existing transaction logs and stored data about the user to generate activity-based authentication questions, our system does not create new vulnerabilities that would affect user privacy.

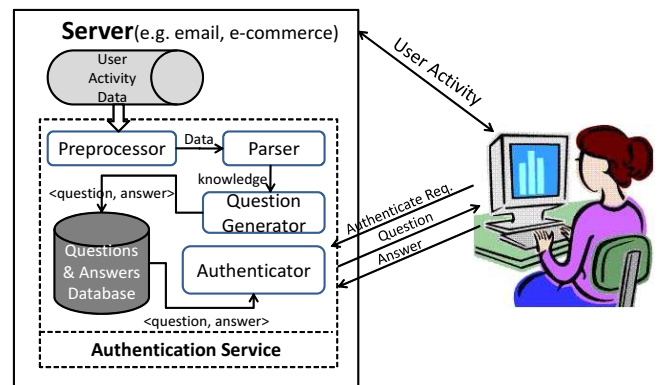


Figure 1: The architecture of an activity-based authentication system

### 4. PRELIMINARY EXPERIMENTS AND RESULTS

In order to study properties of our activity based authentication questions, we conducted a survey with four participants, each knowing each other to some degree through a common advisor. The questions found on the survey can be located in Table 1. The study evaluated activity-based questions on their robustness to attacks and memorability.

Each participant was asked to open a survey in a file that consisted of 12 questions, four from each activity based question

**Table 1: Average vulnerability and memorability levels obtained from the user study (3 being easy to recall the answer and 1 being the opposite).**

Case Study Questions	Correct Guesses	Average Memorability
<i>Network Activity</i>		
1. What was the last website you visited?	0%	2.25
2. Who was the last person you sent an email to?	42%	2.50
3. At what time did you send out your most recent email?	0%	2.00
4. What website do you visit the most?	58%	2.33
<i>Physical Events</i>		
5. What event do (did?) you have planned for Saturday (or day of an event in your calendar)?	0%	2.25
6. Who are you meeting on Tuesday (or day of an event in your calendar)?	41%	2.50
7. Where were you last Monday?	0%	2.00
8. How long was your last meeting scheduled for?	58%	2.33
<i>Conceptual Opinion</i>		
9. What political party do you support?	25%	2.25
10. What is your preferred online news source?	0%	2.50
11. Who is your favorite reporter/blogger?	0%	1.67
12. What is your opinion on North Korean Nuclear testing?	25%	2.00

category (network activity, physical events, and conceptual opinions). The survey had four columns next to each question in which the participant and the three other participants' names were listed across. The objective was for them to answer each question under the column with their names as truthfully as they could and give a memorability ranking (3 being easy to remember the answer and 1 being low memorability), while in the remaining three columns, they were asked to guess the other participants' answers.

We asked our participants to answer a total of twelve questions shown in Table 1. Of the twelve questions the users were able to give answers for all of them with an average memorability level of 2.23. This result shows that the participants had primarily positive reactions in term of the memorability of the questions being asked of them. A more detailed breakdown of participant memorability rankings can be found in Table 1 along with average successful rates of attacks.

Questions found to be more *temporal-based* where the most robust of our activity based authentication questions. Questions 1, 3, 5, and 8 in table 1 asked about time and location, all of which had relatively low rates of attacks. These personal (physical) events that were not related to the work of the participants and with times that also were not related to their work lives dramatically improved the success of these tests in Table 1. In a practical prototype, the physical events may be extracted from a digital calendar or a GPS-enabled mobile device. With an even larger pool of participants, we expect the rate of successful attacks to be even lower.

We found that in this particular user-group work-related questions were the most vulnerable, as the participants were colleagues. On average, our participants' connection to their advisor caused the security problem on some questions that can be easily guessed. This phenomenon was found to be the case in questions 2, 6, and 8 in Table 1. In all of these, they were able to obtain the correct answer by guessing at random or putting down their common link – their advisor. For question 4 in Table 1, the pervasive popularity of certain websites such as Google or Yahoo makes it susceptible attack. We also found that opinion-based questions were relatively hard to guess correctly among the participants. Our preliminary results yield encouraging results on certain types of activity-based authentication questions and motivate us to carry out more thorough investigations on this topic.

## 5. CONCLUSIONS AND FUTURE WORK

In this paper, we introduced an activity-based authentication framework based and described our preliminary evaluation results on its security and usability. Our approach for improving the robustness of existing question-based authentication systems was to use *short-lived* questions that are automatically extracted from the user's personal Internet activities.

In the future, we plan to compare conventional authentication questions with ours by running another case study that allows for individuals to research (e.g., via Google) possible answers before attacking the other participants using public resources. We also plan to expand our study to more participants with diverse backgrounds (as opposed to just colleagues). We plan to implement a prototype of our system with the integration of semantic Web [2, 4] and natural language processing techniques [3, 5]. This prototype will provide further proof of our conceptual system and show the possibility of having an architecture that runs in the client and within the server as well. Our prototype will start with an email server, so that we can extract questions from the user's email logs and calendars. With popular email providers such as Gmail and Yahoo providing calendars as a part of their services, this design would make gathering temporal and physical information on the user a simple process. With these improvements, we hope to see future success in providing more robust activity based authentication questions. Besides being used for server-side authentication such as Web services, we will also explore the potential application of our solution in building a host-based detection system against malicious botnets. The goal of such a system is to challenge the user with a series of questions that will be used to differentiate the legitimate human-user from an invisible bot intruder.

## 6. REFERENCES

- [1] E. Breck, Y. Choi, and C. Cardie. Identifying Expressions of Opinion in Context. In *Proceedings of Twentieth International Joint Conference on Artificial Intelligence (IJCAI)*, 2007.
- [2] D. Fensel, F. van Harmelen, I. Horrocks, D. L. McGuinness, and P. F. Patel-Schneider. Oil: An Ontology Infrastructure for The Semantic Web. *Intelligent Systems, IEEE [see also IEEE Intelligent Systems and Their Applications]*, 16(2):38–45, 2001.
- [3] D. Jurafsky and J. H. Martin. *Speech and Language*

*Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition (Prentice Hall Series in Artificial Intelligence)*. Prentice Hall, 1 edition, February 2000.

- [4] A. Maedche and V. Zacharias. Clustering Ontology-Based Metadata in the Semantic Web. In *PKDD '02: Proceedings of the 6th European Conference on Principles of Data Mining and Knowledge Discovery*, pages 348–360, London, UK, 2002. Springer-Verlag.
- [5] C. D. Manning and H. Schütze. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA, USA, 1999.
- [6] OpenNLP <http://opennlp.sourceforge.net/> 2008.
- [7] S. Schechter, A. J. B. Brush, and S. Egelman. It's No Secret. Measuring the Security and Reliability of Authentication via Secret Questions. In *Proceedings of IEEE Symposium on Security and Privacy*, pages 375–390, 2009.
- [8] A. Whitten and J. Tygar. Why Johnny Can't Encrypt: a Usability Evaluation of PGP 5.0. In *8th Usenix security symposium*, pages 169–184, 1999.