# Opinions or Algorithms: An Investigation of Trust in People Versus Automation in App Store Security

David Schuster[1(✉)], Mary L. Still[1], Jeremiah D. Still[1], Ji Jung Lim[2], Cary S. Feria[1], and Christian P. Rohrer[2]

[1] San José State University, San Jose, CA, USA
{david.schuster,mary.still,jeremiah.still}@sjsu.edu
[2] Intel Security, Santa Clara, CA, USA
diannejlim@gmail.com, christian.p.rohrer@intel.com

**Abstract.** Mobile application (app) stores are a critical source of information about risk in an uncertain environment. App stores ought to assess and communicate the risk associated with an installation so that users are discouraged from installing risky or harmful apps in app stores. However, only a limited number of studies offer designers information about how to communicate risk effectively. We focused on the user's trust associated with security information stemming from crowd-sourced evaluations compared to those generated from an automated system. Both of these sources of security information are pervasively used to indicate possible risk associated with an app. We investigated whether biases exist for a particular source of information given similar amount of security information being available. We found that participants preferred to install apps rated by automation to those rated by humans despite equivalence in stated risk. Further, we found evidence of a gender difference in trust in automation.

**Keywords:** Mobile device security · App stores · Trust in automation · Interpersonal trust

## 1 Introduction

Mobile application (app) usage has become ubiquitous in society, therefore it is increasingly important for users to be able to identity those that pose a security threat. Increases in the number [1] and capabilities of mobile devices have lead to a proliferation of app stores. These stores provide a centralized source for discovering, purchasing, and installing apps [2]. While the security of unknown applications has been long been a concern for desktop computer users, fears could be assuaged by selecting from established brands purchased from brick-and-mortar stores or speaking with those more knowledgeable. Now, app stores provide users with a large number of unknown applications from unfamiliar brands. Further, it may be unclear whether an app's business model is based on gathering personal information or defrauding users. Users must depend on the app stores to protect them from malicious software and to clearly communicate possible risk.

It is not surprising that as the number of available mobile apps increases, so does the prevalence of malware [3, 4]. It appears users are among the last lines of defense in their

own mobile security. They must rely upon information gathered through interactions with app stores to make decisions about the security implications of the apps they download. These interactions with the app store are critical for successful mobile device security.

Unfortunately, app stores may not provide effective and usable communication about the associated installation risks. For example, in their test of apps from the Google Play store, Felt, Chin, Hanna, Song and Wagner [3] found that approximately one-third of Android apps are over-privileged. That is, they present a request for unnecessary permissions (e.g., location data for a flashlight app). Nevertheless, research has shown that users do assume apps are safe and dismiss security warnings [5].

In addition, the typical users of app stores are not security experts, and therefore, most do not know how to interpret security information [6]. The information provided to support their decision to install an app needs to be jargon free and needs to transparently communicate importance. Lin and colleagues [7] suggested, "users have very little support in making good trust decisions regarding what apps to install" (p. 501). It appears they are often making app installation decisions simply based on perceived usefulness [5]. Further, less than 11 % of users in Mylonas and colleagues' study considered the provided reviews, reputation, or security associated with the app. The users who do express concern about security depend on app store community reviews and ratings beyond brand familiarity for establishing trustworthiness before installation. Unfortunately, fellow users, not security experts, generate this content. Despite this, users need a transparent and consistent way to be informed about the variety of possible risk associated with installing a particular app. According to Chia, Yamamoto, and Asokan [8], embedding the Web of Trust service into app store platforms has the potential to better inform users of risk during the app selection decision making process. Web of Trust provides crowd-sourced ratings of web sites as users browse [9]. It appears providing quality information is key to encouraging safe decisions.

In this study, we aimed to inform the presentation of security-relevant information in mobile app stores by comparing two sources of this information: ratings of apps provided by end users and ratings of apps provided by automated methods that evaluate mobile app security.

## 1.1 User-Generated Versus Automated Security Ratings

App security information can be based on user-generated reviews or automated methods. We focus on the distinction between security information generated using a crowd-sourced method and security information generated using an algorithm. In the former source of information, other users of the technology review and rate applications they use. In the latter source, a software agent examines the app and provides a conclusion about its riskiness in terms of security and privacy.

User-generated ratings provide a mechanism for other users to obtain data about other users' experience with an app. Because users share similar goals, the language they use to describe the costs and benefits of using an app may be easier for others to understand. Further, app store owners need only provide a mechanism for users to rate apps; they do not have to develop and maintain their own ratings and reviews. A limitation of this approach is that it requires a large pool of users to provide ratings with sufficient reliability to be of use. Accumulating a large number of user ratings takes time. In the interim,

such as when a given app is new to the store, users must download apps without the benefit of ratings or on the basis of few ratings. The reliability of these ratings is, unfortunately, limited [8].

The benefit of technology-generated ratings, like other forms of automation, is that once an algorithm is developed, providers of an app store can rate apps instantly without delays associated with human input. Technology-generated ratings may work like a virus scanner by looking for suspicious patterns in the app [10] or by comparing the app to known malware. The algorithm may incorporate user feedback and permissions use to assign a trustworthiness metric, as proposed by Kuehnhausen and Frost [11]. In another automated method, software examines the functionality of an app and compares it to the functionally of other known apps [12]. A large difference between the functionality of the app and similar apps could be indicative of over-privilege or malware. Even with these benefits, automated methods share a weakness of user-generated ratings in that they are imperfect. They may miss threats that are present or incorrectly detect threats that are not present (i.e., a false alarm or false positive).

Because both user-generated and technology-generated ratings are used in app stores, and both methods provide imperfect information about the security of apps, an important question is whether users will follow the recommendations of one potentially inaccurate method or another when the warning presented by each is the same. At issue for developers and users is whether people differentially trust people or computers to provide security information.

## 1.2   Trust in App Selection

Trust in an app store provider has been shown to be a factor in user decisions to install apps [13], but within an app store, questions remain about how users trust security ratings to make decisions about which apps to install.

Trust is "the attitude that an agent will help achieve an individual's goals in a situation characterized by uncertainty and vulnerability" (p. 51) [14]. Trust is adaptive; it allows us to make decisions and accept risk under conditions of uncertainty [15] and it allows us to benefit from the effort and expertise of others [14]. Using the term *agent* broadly, several forms of trust are relevant in app store ratings.

The first of these, interpersonal trust, is trust amongst humans [16]. Interpersonal trust describes an expectancy of one human agent that the communications from another human agent are reliable. Although humans are the source of the information, user-generated security ratings are a form of technology-mediated communication. It is important to note differences between trust in a human providing a rating and trust in a medium used to communicate messages from a human. Patrick [15] suggested that both of these are important, and that trust is more difficult to establish when communication is mediated by technology, which he called a once-removed transaction. This is distinguished from non-removed transactions, which involve direct communication between two people. User-generated ratings within an app store are thus an example of a once-removed transaction where interpersonal trust may be a factor in the decision to install an app.

A second form of trust is one of human trust in technology, or trust in automation. It has been repeatedly shown that people base their trust in automation in large part based on perceptions of the performance of the technology [17, 18]. In general, the reliability of

automated aids has a strong affect on people's ability to use it [19]. That is, people tend to trust and use technology more if they perceive it to be reliable [20]. Unfortunately, attributions of technology reliability are not always accurate, especially when users have a limited number of interactions with the technology. This can lead to mis-calibrated trust [21]. When trust is mis-calibrated, people may overtrust, leading to use of a technology beyond its capability, or they may distrust and disuse technology that could be helpful to them [14].

Interpersonal trust (including trust in once-removed transactions) and trust in automation share some similarities (see [14]), but because of important differences between humans and machines, interpersonal trust is affected by different factors than trust in automation [22]. For example, trust in people may be less constrained than trust in automation; a person might be unconditionally trusted across a wide variety of scenarios, but machines are only trusted to do certain tasks [22]. Lewandowsky, Mundy, and Tan [23] found that people felt more responsibility when they believed they were working with automation instead of another person in a simulated pasteurization plant. This suggests that trust is more important for delegation to automation than to a human.

In addition to the qualities of the task and environment, individual differences also affect trust. One such difference is a tendency to trust, a stable trait that also predicts a person's ability to properly calibrate trust [14]. The literature distinguishes between the tendency to trust in other people [16] and the tendency to trust machines [17]. These traits are distinguishable from trust in a particular person, technology, or setting, which may be fluid and change with experience [14].

Gender is another individual difference examined in this literature. Although there is a limited theoretical basis for gender effects in automation trust, gender is commonly collected demographic information that could influence the degree to which a rating method is useful. A study on technology acceptance found that women perceived e-mail as easier to use and more useful than men, but that these differences did not significantly affect e-mail use [25].

In all, the literature suggests that interpersonal trust and trust in automation are separate constructs. Further, we can distinguish between trust in a particular person or tool from a tendency to trust. However, research has not examined whether user-generated or automated methods are trusted more, or used more, at similar levels of uncertainty and risk. To examine this process, we presented users with a series of apps in a simulated app store. Our primary question was whether, when given equivalent alternatives, users would select apps with user-generated ratings or apps with automated ratings. A secondary question was which trust constructs and individual differences predict a decision to rely on user-generated or automated ratings.

## 2    Method

### 2.1    Participants

Forty undergraduate students aged 18–22 years (19 female and 21 male; aged 18–22 years; $M = 19$, $SD = 1.27$; two did not report age) participated in the experiment in exchange for course credit. No participants reported having color-deficient vision.

All participants reported owning at least one device (e.g., tablet, phone) that runs apps. Participants reported downloading an average of two apps per month ($M = 2.17$, $SD = 1.70$).

## 2.2 Materials and Procedure

Participants in this study completed several measures of trust and were presented with simulated situations in which they selected apps in the context of potential risk. First, participants completed a measure of trust in automation [24], a measure of interpersonal trust [16], and a series of demographic questions.

After completing these initial surveys, participants were presented with a simulated app store; this was done on a desktop computer. Participants were given the following instructions, "In this task, you are selecting apps to install on your mobile device. However, apps might have security risks. You will be shown a series of apps, four at a time. For each set of apps, select the safest one. You should always select one of the apps, and you cannot select more than one. To help you decide, security ratings are provided with each app."

Participants were asked to download one app in each trial. During the experiment, participants completed 24 discrete trials. On each trial they were presented with four apps and were to select one of the apps for download. All four apps in a trial belonged to the same category: social media, finance, news, or media player. After making their selection, the next trial was presented.

Of the 24 trials, half were critical trials and half were filler trials. Every trial contained two apps that had user-generated ratings and two apps with automated ratings. On critical trials, two apps, or all four of the apps, were tied for the lowest level of risk according to their ratings. Thus, to select the safest app, participants were forced select between a human and an automated rating.

App safety was indicated by a general security rating with redundant color-coding: "Safe" was displayed in a green font, "Caution" was displayed in orange font, and "Risky" was displayed in red font. Immediately following the security ratings was a verbal description of the threat (see Figs. 1 and 2). Four general types of security threats were used: presence of malware, ability to access and modify account information, perceived breech of privacy, and ability to access and modify phone states (e.g., location services). The apps within one trial had matching threat types (e.g., all might present a perceived privacy breech). An icon and label were used to indicate which ratings were from humans and which were automated. To further differentiate the ratings and preserve a naturalistic element in the study, "human" descriptions were gathered from publically available online comments about existing apps. Automated descriptions reflected generic classes of threats that could be detected using algorithms.

Several measures were taken to reduce the chance of strategic effects in participants and to control the influence of other factors. To reduce strategic effects, 12 filler trials were distributed throughout the 12 critical trials. Those trials contained four apps, as did the critical trials, but the user-generated and automated ratings did not have equivalent safety ratings. This helped disguise our use of a forced-choice paradigm. We were also concerned that specific app icons or names might be more appealing to users and
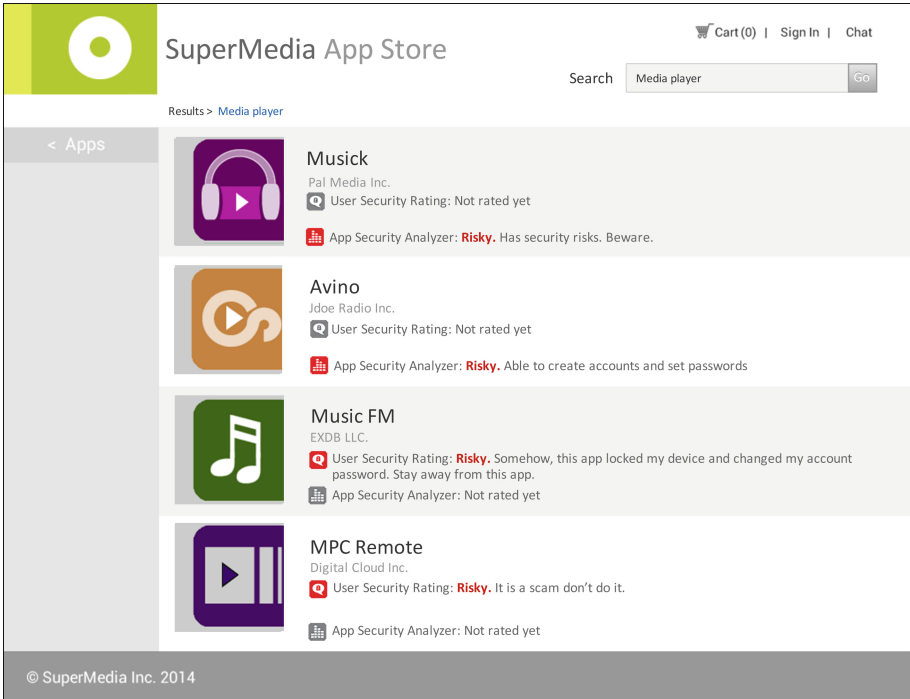
**Fig. 1.** Example of critical trial with equivalent threat levels (*risky*)
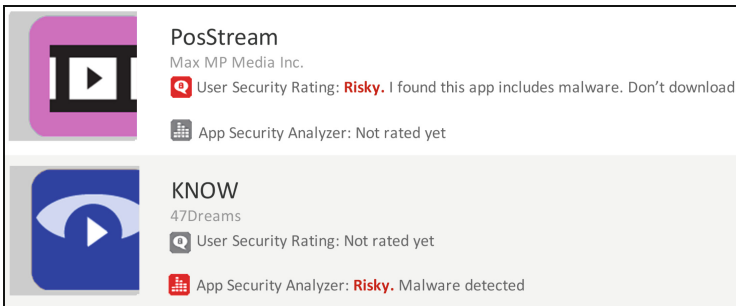


**Fig. 2.** Detail of one *risky* threat rated by the user-generated method (top) and one generated by the automated method (bottom).

confound the data. For example, if most of the preferred icons and names were associated with human-rated apps, it might appear that users trust those ratings more than automated ratings, when in fact, they merely preferred the icon. To control for this, each icon and name was paired an equal number of times with human and automated ratings. This was manipulated between-subjects; for each icon, half of the participants saw it paired with a user rating, and the remaining half of participants saw it paired with an automated rating. Similarly, stimulus presentation was blocked by app category; for instance, all

six social media app trials all six social media app trials were presented one after another. This was done to make the task seem more naturalistic. To control for any order effects created by this blocking technique, we counterbalanced app category order across participants.

After completing the experimental portion of the study, participants completed the checklist for trust between people and automation [26], a self-report survey designed to capture trust in the automated ratings. Participants were instructed as follows: "There are several scales for you to rate intensity of your feeling of trust, or your impression of the computer-generated security ratings in the task you just completed."

## 3   Results

### 3.1   Preference for Apps with the Lowest Risk Rating

Our manipulation of app store rating affected user decisions to select apps as we expected. Participants selected the safest app of the four listed the overwhelming majority of the time ($M = .96$, $SD = .05$). This includes performance on all trials, including distractor trials that were discarded before analyses of rating preference, which we describe next.

### 3.2   Preference for User-Generated Versus Automated Security Ratings

To test our hypothesis that participants would be less willing to trust an automated rating than a human-generated one, we conducted a one-sample t-test to compare the proportion of trials in which participants selected the app with the automated rating to the proportion expected by chance (0.50). This analysis was conducted on critical trials only, those in which the lowest level of security warning was tied between user-generated and automated ratings. When given equivalent human and automated ratings, participants selected apps with automated ratings most of the time ($M = .61$, $SD = .18$), more often than would be expected by chance, $t(39) = 3.77$, $p = .001$.

### 3.3   Individual Differences and Trust

Although participants tended to trust apps rated by automation in comparison to those rated by humans, there may be individual differences in that trust. To explore relationships among trust constructs, we computed a series of Pearson product-moment correlation coefficients. These included the automation induced complacency scale, a dispositional measure of trust in automation [24]. Higher values on this scale indicate greater trust in automation, in general. The checklist for trust between humans and automation [26] measured self-reported trust in the automated ratings used in the study. Higher values on this measure indicate greater trust in a specific technology. Additionally, we included the measure of interpersonal trust [16], the proportion of critical trial automation-rated apps selected (as described previously), and Pearson point-biserial correlations with gender. The results of this analysis are presented in Table 1.

**Table 1.**  Values of Pearson's $r$ and $r_{pb}$ among trust constructs and individual differences

|                        | M     | SD    | 1        | 2    | 3      | 4      |
|------------------------|-------|-------|----------|------|--------|--------|
| 1. Trust in automation | 42.38 | 5.52  | –        |      |        |        |
| 2. Interpersonal trust | 63.43 | 5.02  | −.054    | –    |        |        |
| 3. Checklist for trust | 52.90 | 14.82 | .544***  | .136 | –      |        |
| 4. Automated ratings   | 0.61  | 0.18  | .082     | .178 | −.085  | –      |
| 5. Gender ($r_{pb}$)   |       |       | −.488**  | .090 | −.353* | −.077  |

$N = 40$; * $p < .05$, ** $p < .01$, *** $p < .001$

Trust measures revealed a significant positive relationship between dispositional trust in automation and trust in the automated ratings ($p < .001$). Gender was significantly related to trust; females were less trusting of automation than males ($p = .001$) and less trusting of the automated ratings in the study ($p = .026$).

## 4   Discussion

Previous research suggested users rely on app stores to protect them [5]. Because users remain the last line of defense in their own device security, users need access to information will facilitate better decision-making without requiring expertise in information security (cf. [27]). Trust describes a factor in a users' decision to rely on a source of information given some uncertainty. Our results suggest that users are more likely to trust automated ratings than human-generated ones when both ratings provide seemingly similar levels of risk. This was evident in both users' behavior (i.e., the apps they selected) and in users ratings of trust in the automated rating systems.

The gender differences observed in this study are noteworthy, as females reported less trust in the automated ratings and less trust in automation in general, but we did not observe a similar relationship with preference for automated ratings. Trust is not the only determinant in the decision to use one rating or another, and other factors not measured in this research (e.g., confidence; see [28]) may explain why gender differences in trust did not lead to significant differences in behavior.

Automated ratings are based on imperfect algorithms, and thus they may incorrectly detect threats that are not present or miss threats that are present. However, users may mistakenly believe that the technology-generated ratings are infallible. Thus, if an automated rating states that a particular security problem is present in an app, users may take this as a fact. Over-reliance on automation is a problem of too much trust in automation, a problem that has been demonstrated across domains [14]. At the same time, if a human-generated rating states that a particular security problem is present in an app, users may assume that the human rater could have been mistaken in their assessment of the source of the problem and that the problem might not have been due to the app itself. For instance, if a human rating states that the app changed their account passwords, the user

might question whether the account passwords might have been changed by some other process rather than the app and incorrectly attributed to the app. Also, users might think that the problem described in the rating might have occurred only for that one human rater, and not for other people who have used the app. As people may be more likely to leave comments on crowd-sourced review websites when they have had a negative experience than when their experience was uneventful [29], readers might suspect that the negative occurrence is uncommon.

In this study, each app was displayed with either one user-generated rating or one automated rating. The presence of only one human rating is most representative of a new app for which there has not yet been sufficient time on the market to accumulate numerous ratings. The present results suggest that in the absence of multiple human ratings, users will prefer to base their app security decisions on a technology-generated rating. When only a small number of human ratings are available, users are not able to know how common a problem is or even whether to believe that the problem described in the review is a true problem. However, this lack of trust in human ratings might change if the number of human ratings were to increase, as would be representative of an app that had been on the market for a longer time. Users may scale their trust in user-generated ratings depending on the size of the crowd that contributed to it. Future research should investigate whether users become more likely to base their decisions on human ratings, and have increased trust in human ratings, as the number of human ratings, displayed increases.

# References

1. Statista: Number of apps available in leading app stores as of July 2014 (2014). http://www.statista.com/statistics/276623/number-of-apps-available-in-leading-app-stores/
2. Cramer, H., Rost, M. Bentley, F., Shamma, D.A.: 2nd workshop on research in the large. using app stores, wide distribution channels and big data in UbiComp research. In: UbiComp, pp. 619–620. ACM, New York (2012)
3. Felt, A.P., Chin, E., Hanna, S., Song, D., Wagner, D.: Android permissions demystified. In: 18th ACM Conference on Computer and Communications Security, pp. 627–638. ACM, New York (2011)
4. Zhou, Y., Wang, Z., Zhou W., Jiang, X.: Hey, you, get off of my market: detecting malicious apps in official and alternative Android markets. In: Proceedings of the 19th Network and Distributed System Security Symposium (2012)
5. Mylonas, A., Kastania, A., Gritzalis, D.: Delegate the smartphone user? security awareness in smartphone platforms. Comput. Secur. **34**, 47–66 (2013)
6. Felt, A.P., Ha, E., Egelman, S., Haney, A., Chin, E., Wagner, D.: Android permissions: user attention, comprehension, and behavior. In: Symposium on Usable Privacy and Security, pp. 3–16. ACM, New York (2012)

7. Lin, J., Amini, S., Hong, J.I., Sadeh, N., Lindqvist, J., Zhang, J.: Expectation and purpose: understanding users' mental models of mobile app privacy through crowdsourcing. In: Proceedings of the 2012 ACM Conference on Ubiquitous Computing, pp. 501–510. ACM, New York (2012)

8. Chia, P.H., Yamamoto, Y., Asokan, N.: Is this app safe? a large scale study on application permissions and risk signals. In: Proceedings of the 21st International Conference on World Wide Web, pp. 311–320 (2012)

9. Web of Trust. https://www.mywot.com/en/aboutus

10. Gilbert, P., Chun, B.-G., Cox, L.P., Jung, J.: Vision: automated security validation of mobile apps at app markets. In: 10th International Workshop on Multiple Classifier Systems, pp. 21–26. ACM, New York (2011)

11. Kuehnhausen, M., Frost, V.S.: Trusting smartphone apps? to install or not to install, that is the question. In: IEEE International Multi-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support, pp. 30-37. IEEE (2013) doi:10.1109/CogSIMA.2013.6523820

12. Sarma, B., Li, N., Gates, C., Potharaju, R., Nita-Rotaru, C., Molloy, I.: Android permissions: a perspective combining risks and benefits. In: Symposium on Access control Models and Technologies, pp. 13–22. ACM, New York (2012)

13. Eling, N., Krasnova, H., Widjaja, T., Buxmann, P.: Will you accept an app? empirical investigation of the decisional calculus behind the adoption of applications on Facebook. In: The 34th International Conference on Information Systems. Association for Information Systems (2013)

14. Lee, J.D., See, K.A.: Trust in automation: designing for appropriate reliance. In: Human Factors, vol. 46, pp. 50–80. HFES, Santa Monica (2004)

15. Patrick, A.: Privacy, trust, agents & users: a review of human-factors issues associated with building trustworthy software agents. Technical report, National Research Council Canada (2002)

16. Rotter, J.B.: A new scale for the measurement of interpersonal trust. J. Pers. **35**, 651–665 (1967)

17. Hancock, P.A., Billings, D.R., Schaefer, K.E., Chen, J.Y., De Visser, E.J., Parasuraman, R.: A meta-analysis of factors affecting trust in human-robot interaction. Hum. Factors: J. Hum. Factors Ergon. Soc. **53**, 517–527 (2011)

18. Muir, B.M., Moray, A.N.: Experimental studies of trust and human intervention in a process control simulation. Ergonomics **39**, 429–460 (1996)

19. Johnson, R.C., Saboe, K.N., Prewett, M.S., Coovert, M.D., Elliott, L.R.: Autonomy and automation reliability in human-robot interaction: a qualitative review. In: Proceedings of the Human Factors and Ergonomics Society Annual Meeting, vol. 53, pp. 1398–1402. Human Factors and Ergonomics Society, Santa Monica, CA (2009)

20. Madhavan, P., Wiegmann, D.A.: Similarities and differences between human–human and human–automation trust: an integrative review. Theor. Issues Ergon. Sci. **8**, 277–301 (2007)

21. Dzindolet, M.T., Peterson, S.A., Pomranky, R.A., Pierce, L.G., Beck, H.P.: The role of trust in automation reliance. Int. J. Hum Comput Stud. **58**, 697–718 (2003)

22. Hoffman, R.R., Bradshaw, J. M., Ford, K.M., Underbrink, A.: Trust in automation. In: IEEE Intelligent Systems (2013)

23. Lewandowsky, S., Mundy, M., Tan, G.: The dynamics of trust: comparing humans to automation. J. Exp. Psychol. Appl. **6**, 104 (2000)

24. Singh, I.L., Molloy, R., Parasuraman, R.: Development and validation of a scale of automation-induced "Complacency". In: Proceedings of the Human Factors and Ergonomics Society Annual Meeting, vol. 36, pp. 22–25. SAGE Publications (1992)

25. Gefen, D., Straub, D.: Gender difference in the perception and use of e-mail: an extension to the technology acceptance model. MIS Q. **21**, 389–400 (1997)
26. Jian, J.-Y., Bisantz, A.M., Drury, C.G.: Foundations for an empirically determined scale of trust in automated systems. J. Cogn. Ergon. **4**, 53–71 (2000)
27. Tam, J., Reeder, R.W., Schechter, S.: I'm allowing what? disclosing the authority applications demand of users as a condition of installation. Technical Report, Microsoft Research (2010). http://research.microsoft.com/apps/pubs/default.aspx?id=131517
28. Parasuraman, R., Riley, V.: Humans and automation: use, misuse, disuse, abuse. Hum. Factors **39**, 230–253 (1997)
29. Hu, N., Pavlou, P.A., Zhang, J.: Can online reviews reveal a product's true quality?: empirical findings and analytical modeling of online word-of-mouth communication. In: Proceedings of the 7th ACM Conference On Electronic Commerce, pp. 324–330. ACM, New York (2006)