

Combating Phishing Attacks: A Knowledge Management Approach

Matthew L. Jensen
University of Oklahoma
mjensen@ou.edu

Alexandra Durcikova
University of Oklahoma
alex@ou.edu

Ryan T. Wright
University of Virginia
ryantwright@gmail.com

Abstract

This paper explores how an organization can utilize its employees to combat phishing attacks collectively through coordinating their activities to create a human firewall. We utilize knowledge management research on knowledge sharing to guide the design of an experiment that explores a central reporting and dissemination platform for phishing attacks. The 2x2 experiment tests the effects of public attribution (to the first person reporting a phishing message) and validation (by the security team) of phishing messages on reporting motivation and accuracy. Results demonstrate that knowledge management techniques are transferable to organizational security and that knowledge management can benefit from insights gained from combating phishing. Specifically, we highlight the need to both publicly acknowledge the contribution to a knowledge management system and provide validation of the contribution. As we saw in our experiment, doing only one or the other does not improve outcomes for correct phishing reports (hits).

1. Introduction

Employees all world over receive them: messages that tempt us to click on a link to address emergencies such as preserving our email accounts from deletion or viewing a critical security notification. Some phishing messages are easy to spot, but many deceive even the most skilled computer users [1]. Organizations of all types (commercial, governmental, and nonprofit) are under constant threat from others who wish to steal private information. Increasingly the most dangerous threat for a data breach comes from phishing attacks through legitimate channels of electronic communication such as email [2]. The damages from these attacks, which include fraud, theft, damage to reputation, regulatory violations, and loss of intellectual property [3], have been estimated to exceed \$2.3 billion USD annually according to the FBI [4]. Gone are the days when organizations might successfully avoid attack by keeping a low profile and

maintaining firewalls and servers. The issue facing most organizations is no longer if, but when a phishing attack will reach organization members.

Research so far has examined how individuals can avoid these attacks through implementing training and security policies [e.g., 5] and SANS institute labeled this initiative as the ‘human being firewall’ [6]. Results from these studies point to a reduction in phishing vulnerability with training and security policies. However, research also suggested that even with training, a few people when left alone still struggle with identifying phishing attacks [e.g., 7]. For a phishing attack to be successful, often all that is required is for a single person in a targeted group to respond.

This research investigates a different approach in which *individuals work together*, rather than in isolation as suggested by the SANS institute initiative [6], to create an interconnected “human firewall”. Individuals acting together directly addresses the problem of the weakest link by which phishers often achieve success. With organization members acting in a coordinated fashion and sharing knowledge about attacks, an individual need not face phishing attacks alone, but can be informed and protected by other organization members.

The creation of the human firewall builds on previous research that shows people can recognize phishing attacks quickly when acting as a group [8]. However, how identification and dissemination can be facilitated by technology is an unexplored issue that needs to be resolved if organizations can make use of a human firewall. Therefore, we draw upon strategies from knowledge sharing and information security [e.g., 9] to guide our investigation. By blending these two perspectives in novel ways, our hope is to improve coordination between organization members as they face evolving phishing attacks.

In this research, we answer the following research question: *How can organizations leverage knowledge sharing technologies and extrinsic motivation to facilitate the human firewall by sustaining organization members’ motivation to contribute, increasing correct identification of phishing messages, and reducing incorrect identification of*

phishing messages? To answer our research question, we draw on theory from knowledge management and crowdsourcing to guide our hypotheses. We carried out a 2 x 2 experiment that crossed attribution (to the first person reporting a phishing message) and validation by the security team of phishing messages received by organization members. In answering this research question with our experiment, this research makes the following contributions: First, we test the feasibility of shifting the focus of anti-phishing efforts from individuals to groups of individuals and potentially whole organizations. Second, our research tests technology-based, organizational interventions (e.g., attribution and validation) that managers may consider to facilitate individuals working together to counter phishing attacks.

2. Background Literature

2.1. Vulnerability to Phishing Attacks

Phishing is “a form of social engineering in which an attacker, also known as a ‘phisher’, attempts to fraudulently retrieve legitimate users’ confidential or sensitive credentials by mimicking electronic communications from a trustworthy or public organization in an automated fashion” [10, p.1]. The Anti-Phishing Working Group (APWG) reported that they observed more phishing attacks in the first quarter of 2016 than in any other three month period since 2004 when they started collecting data. In addition, APWG reported the number of unique phishing websites detected increased 250 percent between October 2015 and March 2016 [11]. Ponemon Institute estimates annual cost of a successful phishing attack per company to be \$3.7 million while about 48% of this cost relates to costs from loss of employee productivity [12]. Spear phishing attacks that have been successful at companies and institutions such as Target, Sony, and even the Pentagon and White House, cost on average around \$1.8 million per incident [13]. Despite the tremendous rate of growth, Vishwanath and colleagues [14] point out that the prevalence of phishing attacks diminish consumer confidence and trust in online commerce and communication, resulting in increased operational costs for online retailers. Thus, research that focuses on how to combat these types of attacks is a top priority not only for researchers but also for IT practitioners.

2.2 Gathering Phishing Knowledge

Researchers have been studying what motivates employees to submit knowledge to a central repository for decades. Most of this research has come through investigation of knowledge management systems (KMSs) because without new inflow of knowledge a KMS cannot deliver value [15]. This is especially true for KMSs that support a fast changing environment, such as tracking phishing attacks. There are two required steps to getting high-quality knowledge into a KMS: (1) knowledge must be contributed by the members of the community that utilize this knowledge; and (2) the contributed knowledge must be validated to ensure accuracy [16]. Regarding the contribution by organizational members, research found that the following factors play a role: intrinsic and extrinsic motivation [e.g., 17], attitude toward knowledge sharing, subjective norm of knowledge sharing [e.g., 18], individual motivations, structural capital, cognitive capital, and relational capital [19]. The second step, validation, must be carefully designed because if even highly motivated employees perceive the validation process to be too strict or non-transparent, they will stop submitting knowledge [20] because rejection may lead to embarrassment [21] or could be perceived to be very costly [18].

Based on the above, we decided to focus on an important motivational factor – extrinsic motivation [e.g., 22]. Specifically, we focus on attribution of the contribution of potential phishing messages to the first person who reported it. We also incorporated a transparent validation process that provided timely feedback and did not reject any submissions [20] so as to encourage reporting of potential phishing attacks [18].

2.3 Accuracy of Reporting Phishing

The determination of phishing/non-phishing is a binary identification task for which accuracy can easily be measured in terms of hits and false positives. In phishing identification, *hits* occur when individuals correctly classify actual phishing messages as phishing. *False positives* occur when individuals classify legitimate messages as phishing messages. From these two measurements, the other potential outcomes of identification tasks can be easily calculated (e.g., false negative, correct rejection).

Past theorizing regarding identification tasks [e.g., 23] has suggested that there are two primary mechanisms available to individuals who wish to improve accuracy in identification tasks. The first mechanism is to properly set thresholds that individuals use in their identification tasks. While deciding whether or not a message is phishing, an individual may observe several characteristics and if a

combination of these characteristics exceeds thresholds, the individual will classify the message as phishing. For example, an individual could examine such characteristics as source familiarity, language that induces time pressure, and inclusion of a suspicious link. If the presence or amount of these characteristics exceeds a combination of thresholds determined by the individual, then the message would be classified as phishing. Careful placement of the thresholds will improve accuracy, especially if biases or habit may cloud the identification task.

The second way to improve accuracy is to increase the number of diagnostic characteristics that may be useful in the identification task. For example, if the individuals may learn that a request for private information is a highly suspicious request. Therefore, individuals may incorporate the type of request as a useful characteristic to which they should pay attention.

The object of most training programs designed to improve individuals' accuracy in identification tasks is improving the placement of thresholds and introducing new, diagnostic characteristics that should be considered during identification. However, in order for individuals to gauge their performance and internalize lessons, some kind of validation is necessary. When validation is provided, individuals have the opportunity to adjust the characteristics they attend to and the threshold they apply to them. In the phishing context, validation is provided through *validation* as the ground truth regarding a particular message is uncovered and reported.

3. Hypotheses Development

When organizations publically acknowledge contributions from organizational members, they will be incentivized to report more messages to gain public recognition through their interaction with the KMS [24]. Attribution, a type of extrinsic reward [22], communicates to the individual that the organization values their contribution and assigns personal credit. Public attribution is a way to build reputation [19, 25] and is evidence of expertise [26]. As a result, when individuals feel that knowledge sharing can elevate their reputation, they will be more inclined to submit potential phishing attacks to the KMS [e.g., 19].

H1: Public attribution of contribution will increase self-reported motivation to report suspicious messages.

With greater motivation to contribute, individuals will be likely to be more sensitized to the potential for

phishing messages. In other words, the thresholds that individuals use to identify a phishing message may be lowered. The lowered thresholds could result in an increase in the number of hits individuals achieve, but would likely come at a cost of an increased number of false positives.

H2: Public attribution of contribution will increase the number of suspicious emails that are falsely reported as phishing (false positives).

H3: Public attribution will increase the number of correctly reported phishing messages (hits).

Validation of the reported potential phishing attacks serves two purposes that may increase the number of hits and decrease the number of false alarms. First, when others (e.g., IT Security department) review the reported messages and provide validation, individuals will realize that their contributions are being evaluated for correctness. They will become more motivated to carefully process messages they report [e.g., 20]. Second, validation may offer individuals the ability to improve their own decision making as they have a chance to adjust their thresholds and the characteristics they consider. Additionally, when validation is made public, individuals have the chance to not only learn from their own experience reporting messages as suspicious, but they also have the benefit of observing and learning from the validation results of others. Therefore, public validation supports observational and experiential learning, which should lead to an increase in the number of hits and a decrease in the number of false positives.

H4: Public validation will interact with public attribution such that a) false positives will decrease and b) hits will increase when they are combined.

4. Method

To test our hypotheses, we conducted a 2 x 2 experiment, crossing attribution (present or absent) and validation (present or absent). The experiment included two parts. The first was a pre-survey that participants completed several days prior to coming to the lab. The pre-survey contained questions about the covariates and permitted participants to schedule a lab session. The second part was a 40-minute lab session during which participants were asked to assume the role of an intern to a senior vice president (SVP) of a software company. Upon arriving at the experiment, participants were consented and then were introduced

to their tasks. Participants were also given a list of employees and personal contacts for the SVP and each participant was instructed to help manage the SVP's inbox. Participants responded to messages from other executives, scheduled meetings for the SVP, and forwarded personal to the SVP's personal account. In addition participants were instructed to help plan a future product marketing event by finding three different hotels in a remote city that have sufficient capacity to handle the event. The messages in each inbox and all work tasks were the same for all participants. These work tasks were meant simulate the multiple organizational priorities (in addition to information security) that employees must manage.

In addition to their work tasks, participants were asked to read an organizational security policy, which required them to report phishing messages by forwarding them to the IT security department. Participants were instructed that completing the work tasks was a higher priority than reporting suspicious messages.

When participants opened the inbox, there were 8 emails waiting to be processed, one of which was a phishing message. An additional 18 emails were sent to participants including four additional phishing messages. Phishing emails were modeled after actual malicious messages [27] and mimicked an IT-service desk request, a cloud storage share request, a deal from a hotel chain, a payment receipt, and a security alert. All phishing emails contained links to a website owned by the experimenters. If participants clicked on a link in a phishing email, they were first directed to a webpage in our website (where their machines could have been compromised if the phishing attack were real) and then were immediately redirected to a legitimate website. Participants had a total of 30 minutes to process all 26 messages, after which they were directed to a post-survey where they completed items concerning motivation. Finally, participants were asked not to share details of the experiment with others and were dismissed.

4. 1. Participants

Students from an introductory MIS class at a large U.S. mid-western university were recruited for the experiment and were offered extra credit for their participation. A total of 120 students completed the pre-survey however, 16 students did not complete the lab session and were excluded from the study. Therefore a total of 104 completed the experiment by attending the lab session. Participants who completed the experiment reported a mean age of 20.6 (max: 33, min: 18) and of all participants, 72.1 percent were male. Students were a good population for this

experiment because a majority of them would shortly join the workforce as interns and would use email during their work. Additionally, students are a frequent target of phishing attacks. 33.7 percent of participants reported knowing someone who had fallen for a phishing message and 34.6 percent of participants reported nearly falling for a phishing message themselves.

4.2. Independent, Dependent, and Covariate Variables

4.2.1. Independent Variables. The security policy provided to the participants described the experiment manipulations and introduced a message board, acting as a KMS, that all participants in a session could see. The message board, displayed the following information about messages that were reported to IT security: (1) subject of the message; (2) number of people who reported the message; (3) first participant to report the message (in the condition where attribution was present); (4) validation status – “under review,” “confirmed phishing,” “confirmed spam,” “non-malicious” (in the condition where validation was present). In the validation condition, messages that were reported were initially labeled as under review. Ninety seconds after the email was reported, the status changed to the either confirmed phishing, confirmed spam, or non-malicious. To ensure all participants understood the purpose and function of the message board, during the introduction of the experiment an experimenter would present the message board, describe all of its components, and answer any questions participants had. An example message board with all four components is shown in Figure 1.

Figure 1. Screenshot of message board with all conditions shown.

Message/Subject	First Reported By	Number of Reports	Status
News Today	Intern1005	12	Verified SPAM
Security Alert	Intern1013	11	Verified Phishing Message
Industry New URL	Intern1012	10	Verified SPAM
Your co-worker shared a folder with you on Dropbox	Intern1004	10	Verified Phishing Message
New Product Launch	Intern1003	9	Legitimate Email Message
Accident on the Pike	Intern1002	8	Verified SPAM
Running Late	Intern1012	5	Legitimate Email Message
Payment Receipt	Intern1003	5	Verified Phishing Message
IT-service Desk	Intern1013	5	Verified Phishing Message
News Today	Intern1004	5	Verified SPAM
Thank you for the booking our hotel	Intern1015	4	Verified Phishing Message
Tax Info	Intern1003	2	Legitimate Email Message
FYI - Out of town	Intern1009	1	Legitimate Email Message
Beers tonight	Intern1002	1	Legitimate Email Message
Cancel meeting	Intern1009	1	Legitimate Email Message
Filing Meeting	Intern1009	1	Legitimate Email Message
Recommend a Marketing Firm	Intern1009	1	Legitimate Email Message

4.2.2. Dependent Variables. The study addressed three dependent variables. The first is self-reported motivation to report phishing messages. This scale included three items: “I tried hard to identify phishing messages during the task,” “I was motivated to report phishing messages,” and “Reporting phishing messages was important to me.”

The second and third dependent variables captured the hits and false positives from participants as they reported suspicious messages. If a participant reported a phishing message it, was recorded as a hit. If a participant reported a spam message or non-malicious message, it was recorded as a false positive. Repeated reports of phishing and non-phishing messages by the same participant were ignored. Therefore the maximum number of hits a participant could have was five and the maximum number of false positives a participant could have was 21.

4.2.3. Covariates. Past research has demonstrated other factors that influence individuals’ recognition of and resistance to phishing messages [7, 28]. Therefore, we captured propensity to trust [29], perceived Internet risk [30, 31], internal and external computer self-efficacy [32], and self-reported expertise in identifying phishing messages as these variables have been examined in recent phishing research [e.g., 7].

5. Data Analysis and Results

Using Mplus 7.1 [33], a measurement model was estimated to determine reliability, discriminant validity, convergent validity, and calculate latent factor scores for self-reported measures. The fit statistics for the measurement model were CFI = 0.951, TLI = 0.940, RMSEA was 0.056 with a 90% confidence interval of 0.034-0.075. All of these fit statistics provide evidence of convergent validity [34]. Further, construct cross-loadings were analyzed to provide evidence of discriminant validity. All of the loadings of each item on its latent construct exceeded 0.6. Average variance extracted for all constructs was much larger than 0.5; therefore good convergent validity was demonstrated [35], and all square roots of average variance extracted exceeded the correlation coefficients between construct and therefore demonstrated good discriminant validity [36].

The analysis plan consisted of 3 different Analyses of Covariance (ANCOVAs). In each ANCOVA, the latent factor scores were used to estimate perceptual measures. Further, attribution and validation served as the independent variables and propensity to trust, Internet risk, internal computer self-efficacy, external self-efficacy, and expertise in identifying phishing

served as covariates. In the first ANCOVA, motivation served as the dependent variable. In the second, number of hits served as the dependent variable. In the third, number of false positive served as the dependent variable. The means and standard deviations for motivation, hits, and false positives for all experimental conditions are shown in Table 1.

Consistent with H1, the first ANCOVA revealed a significant effect of attribution on motivation, $F(1, 95) = 5.210, p = .003, \eta_p^2 = .09$. However, the second ANCOVA did not reveal a significant effect of attribution on hits, $F(1, 95) = .095, p = .759$, and the third ANCOVA did not reveal a significant effect of attribution on false positives, $F(1, 95) = 2.307, p = .132$. These finding fail to confirm H2 and H3.

Table 1. Means of dependent variables by condition

Condition	N	Mean Motivation (SD)	Mean Hits (SD)	Mean False Positives (SD)
No Attribution, No Validation	25	-.210 (.848)	3.480 (1.123)	2.120 (1.943)
Attribution, No Validation	30	.148 (.599)	2.533 (1.525)	2.300 (1.705)
No Attribution, Validation	25	-.255 (1.064)	2.680 (1.887)	2.040 (2.131)
Attribution, Validation	24	.299 (.407)	3.542 (1.414)	2.958 (1.517)

In the first ANCOVA, external computer self-efficacy exerted a significant influence on motivation, $F(1, 95) = 3.594, p = .014, \eta_p^2 = .06$. But all other covariates were insignificant.

In the second ANCOVA, internal computer self-efficacy exerted a significant influence on hits, $F(1, 95) = 4.426, p = .038, \eta_p^2 = .05$, and the influence of external computer self-efficacy approached significance, $F(1, 95) = 3.172, p = .078, \eta_p^2 = .03$. All other covariates were insignificant.

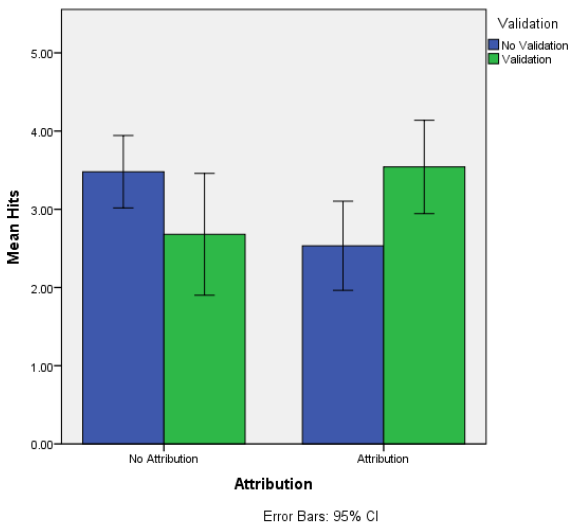
In the third ANCOVA, all of the covariates were insignificant.

To test the H4, we examined the interaction coefficient for attribution and valuation in the second and third ANCOVAs. The influence of the interaction coefficient on false positives was insignificant in the third ANCOVA, $F(1, 95) = 1.215, p = .273$. However, the influence of the interaction coefficient on hits was significant, $F(1, 95) = 8.027, p = .006, \eta_p^2 = .08$. The interaction is shown in Figure 2.

Simple pairwise comparisons show that participants using a message board without attribution

and validation had more hits than those in the attribution-only and validation-only conditions (see Table 1 for means). Likewise, participants in the attribution and validation condition had more hits in than the attribution-only and validation-only conditions. Interestingly, participants using a message board with neither attribution nor validation had a similar number of hits as participants with both.

Figure 2. Interaction of attribution and validation on number of hits



5.1. Supplemental Analysis

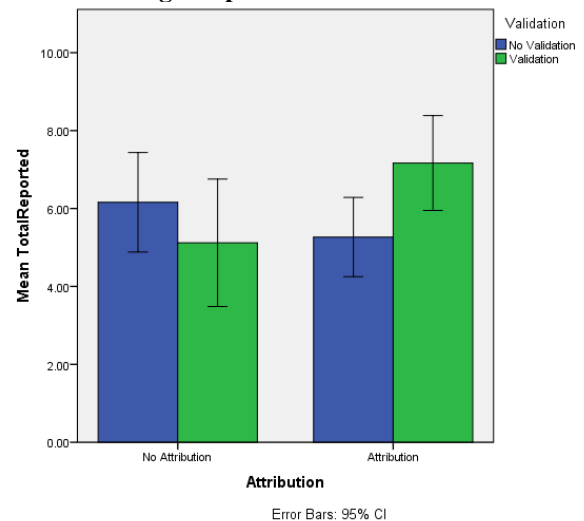
To explore further the effects of attribution and validation, we examined their effects on the level of participation (whether or not participants reported any emails as suspicious) and the raw counts of messages that were reported as suspicious (taking into account duplicate reports).

To examine participation, we conducted a logistic regression with whether or not each participant reported a message as suspicious as the dependent variable and attribution and validation as independent variables. All of the covariates were also included. Results showed that 100 out of the 104 participants participated by reporting at least a single message as suspicious. Not surprisingly, neither attribution ($B = -19.40, p = .998$) nor validation ($B = -17.84, p = .998$) produced a significant effect on participation. None of the covariates were significant either.

To examine the raw number of reports, we conducted an additional ANCOVA with attribution and validation as independent variables. We also include the covariates in the analysis. Results revealed no significant main effects or significant covariates.

However, the interaction between attribution and validation was significant, $F(1, 95) = 4.972, p = .028, \eta_p^2 = .05$. The interaction followed a similar pattern to the interaction pattern produced in the ANCOVA for hits. Simple comparison tests revealed that participants in the condition with both attribution and validation reported more messages as suspicious as those in the attribution-only condition and those in the validation-only condition. However, there was no difference between participants with both attribution and validation and participants with neither. The interaction is illustrated in Figure 3.

Figure 3. Interaction of attribution and validation on total messages reported

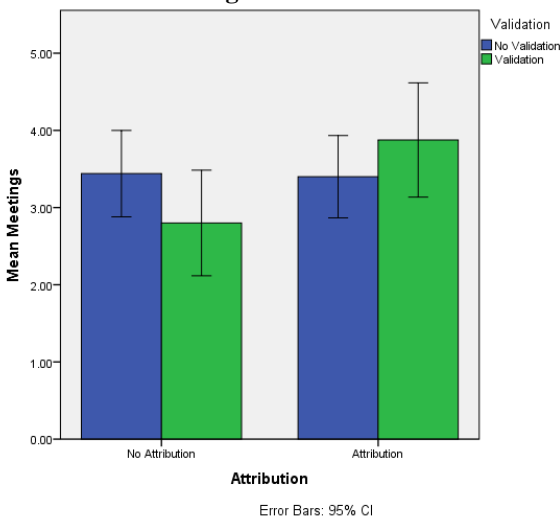


Finally, to determine if conditions of the message board were disruptive to work tasks, we conducted additional analysis on the number of messages the participants sent to co-workers and the number of meetings scheduled as part of their work tasks. We did not find significant differences in the number of messages the participants sent due to attribution, $F(1, 95) = .027, p = .869$, or due to validation, $F(1, 95) = 1.015, p = .316$. Similarly, we did not find significant differences in the number of meetings participants scheduled due to validation, $F(1, 95) = .001, p = .982$. However, the effect from attribution approached significance, $F(1, 95) = 2.905, p = .092, \eta_p^2 = .03$, and the attribution x validation interaction was significant, $F(1, 95) = 3.894, p = .051, \eta_p^2 = .04$. The interaction is illustrated in Figure 4. Simple comparisons demonstrated that those in the attribution, validation condition scheduled significantly more meetings than those in the no attribution, validation condition.

6. Discussion

Before discussing the implications of our findings for research and practice, we raise several limitations of this study. First, our study is subject to many of the limitations common to experimental research. The participants were not actually employed by an organization and were playing a role. Although the role would have been familiar to participants, they were not subject to many of the organizational pressures that actual interns would face. These limitations permitted assignment to conditions and enabled experimental control.

Figure 4. Interaction of attribution and validation on number of meetings scheduled



Second, we anticipate that the results of this study generalize to the college-aged student population. However, additional work is necessary to determine if these results will generalize to a more diverse sample of working adults.

Third, our focus was on regular phishing attack and not spear phishing attacks that are more successful. In our study, the success of these regular phishing attacks was 25%, which means that out of 5 phishing emails subjects fell for at least one. Future research should also focus on spear phishing attacks.

With these limitations in mind, there are several important implications of this study.

6.1 Implications for Research

We have several key findings for both knowledge sharing and security. First, we found that attribution within our message board does positively affect motivation to report phishing messages. Past research

on contribution to KMSs states that “individuals must think that their contribution to others will be worth the effort and that some new value will be created” [37 pg. 36]. Further, the literature states that personal acknowledge and personal benefits do indeed increase motivation to contribute [38]. Our study provides a message board that shows contribution activity, so you can see the collective group participating in real time. This alone creates motivation for participants to report phishing messages.

We did find that the conditions which produced the best hit rates (e.g., identifying of phishing messages) were either no attribution and validation or both attribution and validation. Attribution or validation alone were the least likely elicit hits.

Bock et al. [39] argue that several factors drive attitudes towards knowledge sharing. These include: (1) anticipated rewards, (2) anticipated reciprocity and (3) self-worth. Bock describes self-worth as an employee getting feedback on their contribution will develop a favorable attitude towards contribution. We also know that individuals working on a collective task feel they are central to the effort are more likely to sustain their contribution [40]. Taken together and as also supported by the results of our experiment, you need to acknowledge the contribution both publically and provide validation on the contribution [41]. In addition, as we saw in our experiment, doing only one or the other does not improve outcomes. Interestingly, previous research showed that attribution and validation are peripheral cues that are more important on knowledge filtering decisions than is the content of a knowledge submission [41]. Our research shows that attribution and validation are also important for knowledge submission. This is an interesting finding because it shows that the interaction between attribution and validation influences both knowledge reuse and knowledge sharing, two processes that were deemed by the literature to have completely different antecedents. Future studies need to evaluate whether these two factors apply to all kinds of knowledge (e.g., elaborate documents, short explanations of what to do next) and whether they apply in different knowledge domains (e.g., security, project management, computer help-desk).

We also found that attribution for phishing reporting created spillover to other non-related tasks. Specifically, we found that once activated by attribution, participants increased the completion of work-related tasks. That said, the effect of attribution and validation did not spill over to identifying the other work related tasks. This spillover was an unexpected benefit of the public message board that warrants additional investigation.

In conclusion, our lab experiment in which we developed a leaderboard-like system to track phishing messages that functioned as a knowledge management system that supported the human firewall not only contributes to the security literature that studies how to combat phishing attacks but also to the knowledge sharing literature.

6.2 Implications for Practice

Clearly, there are effective alternatives to incentivize contributions to centralized anti-phishing efforts. This work provides evidence for this. Organizations must take note that providing a message might add value, but alone it is not optimal. Further, organizations must understand the implications of adding certain design features that are guided by current research in knowledge sharing [16-20] and knowledge reuse [41], such as attribution of shared knowledge to the author and validation of such knowledge.

Second, motivation alone does not improve the security of your organizations. Many past studies have argued that motivation and fear appeals are the linchpins to improved IT security [42, 43]. We found manipulations that increase one's motivation, but they did not necessarily improve the phishing message hits or even the false positives. Organizations need to be aware that solutions that may, at face, increase the motivation of your workforce to report vulnerabilities, may not mitigate the vulnerabilities.

7. Conclusion

This preliminary study examines whether certain design elements in a knowledge management system (e.g., message board) are transferable to support the human firewall where individuals work together, rather than in isolation, to combat phishing attacks. We found that attribution of the contribution increases motivation to contribute, but not the overall quality of the contribution. The optimal design in our study was attribution and validation of the contribution or a plain message board with neither of these design elements.

This research was supported by a grant from the NSF Social and Economic Division, Project# 1421580. The views and conclusions contained herein are those of the authors and should not be interpreted as representing the official policies or endorsements, either expressed or implied, of NSF or the U.S. Government.

10. References

- [1] Jackson, C., Simon, D.R., Tan, D.S., and Barth, A., "An Evaluation of Extended Validation and Picture-in-Picture Phishing Attacks": Financial Cryptography and Data Security, 2007, pp. An Evaluation of Extended Validation and Picture-in-Picture Phishing Attacks.
- [2] Jagatic, T., Johnson, N., and Jacobsson, F., "Social Phishing", Communications of the ACM, 50(10), 2007, pp. 94-100.
- [3] Hong, J., "The State of Phishing Attacks", Communications of the ACM, 55(1), 2012, pp. 74-81.
- [4] <https://www.fbi.gov/phoenix/press-releases/2016/fbi-warns-of-dramatic-increase-in-business-e-mail-scams>, accessed June 14, 2016.
- [5] Vance, A., Siponen, M., and Pahlila, S., "Motivating Is Security Compliance: Insights from Habit and Protection Motivation Theory.", Information & Management, 49(3), 2012, pp. 190-198.
- [6] El-Harmeel, M., "Human Being Firewall", in (Editor, 'ed.'^eds.): Book Human Being Firewall, 2009
- [7] Wright, R.T., and Marett, K., "The Influence of Experiential and Dispositional Factors in Phishing: An Empirical Investigation of the Deceived", Journal of Management Information Systems, 27(1), 2010, pp. 273-303.
- [8] Kelley, C.M., Hong, K.W., Mayhorn, C.B., and Murphy-Hill, E., "Something Smells Phishy: Exploring Definitions, Consequences, and Reactions to Phishing", in (Editor, 'ed.'^eds.): Book Something Smells Phishy: Exploring Definitions, Consequences, and Reactions to Phishing, 2012, pp. 2108-2112.
- [9] Jennex, M., and Durcikova, A., "Integrating Km and Security: Are We Doing Enough?" ", International Journal of Knowledge Management, 10(2), 2014, pp. 1-12.
- [10] Myers, S., "Introduction to Phishing", in (Jakobsson, M., and Myers, S., 'eds.): Phishing and Countermeasures: Understanding the Increasing Problem of Electronic Identity Theft, Wiley, Hoboken, NJ, 2007, pp. 1-29.
- [11] Apwg, "Phishing Activity Trends Report: 1st Quarter", in (Editor, 'ed.'^eds.): Book Phishing

Activity Trends Report: 1st Quarter, APWG, www.apwg.org, 2016

[12] Samarati, M., "Cost of Phishing Estimated to Be \$3.7 Million", in (Editor, 'ed.'^eds.): Book Cost of Phishing Estimated to Be \$3.7 Million, 2016

[13] Seals, T., "Average Cost of a Spear Phishing Incident: \$1.6 Mn", in (Editor, 'ed.'^eds.): Book Average Cost of a Spear Phishing Incident: \$1.6 Mn, 2016

[14] Vishwanath, A., Herath, T., Chen, R., Wang, J., and Rao, H.R., "Why Do People Get Phished? Testing Individual Differences in Phishing Vulnerability within an Integrated, Information Processing Model", *Decision Support Systems*, 51(3), 2011, pp. 576-586.

[15] Olivera, F., "Memory Systems in Organizations: An Empirical Investigation of Mechanisms for Knowledge Collection, Storage and Access", *Journal of Management Studies*, 37(6), 2000, pp. 811-832.

[16] Fadel, K., and Durcikova, A., "If It's Fair, I'll Share: The Effect of Knowledge Validation Justice on Contribution to a Knowledge Repository", *Information & Management*, 51(5), 2015, pp. 511-619.

[17] Kankanhalli, A., Tan, B., and Wei, K.K., "Contributing Knowledge to Electronic Knowledge Repositories: An Empirical Investigation", *MIS Quarterly*, 29(1), 2005, pp. 113-143.

[18] Bock, G.W., Kim, Y.G., and Lee, J., "Behavioral Intention Formation in Knowledge Sharing: Examining the Roles of Extrinsic Motivators, Social-Psychological Forces and Organizational Climate", *MIS Quarterly* 29, 1(87-111), 2005,

[19] Wasko, M., and Faraj, S., "Why Should I Share? Examining Social Capital and Knowledge Contribution in Electronic Networks of Practice", *MIS Quarterly*, 29(134-56), 2005,

[20] Durcikova, A., and Gray, P.H., "How Knowledge Validation Processes Affect Knowledge Contribution", *Journal of Management Information Systems*, 25(4), 2009, pp. 81-107.

[21] Husted, K., and Michailova, S., "Diagnosing and Fighting Knowledge-Sharing Hostility", *Organizational Dynamics*, 31(1), 2002, pp. 60-73.

[22] Hung, S.-Y., Durcikova, A., Lai, H.-M., and Lin, W.-M., "The Influence of Intrinsic and Extrinsic Motivation on Individuals' Knowledge Sharing Behavior", *International Journal of Human-Computer Studies*, 69(6), 2011, pp. 415-427.

[23] Jensen, M.L., Lowry, P.B., and Jenkins, J.L., "Effects of Automated and Participative Decision Support in Computer-Aided Credibility Assessment", *Journal of Management Information Systems*, 28(1), 2011, pp. 203-236.

[24] Ba, S., Stallaert, J., and Whinston, A.B., "Research Commentary: Introducing a Third Dimension in Information Systems Design - the Case for Incentive Alignment.", *Information Systems Research*, 2(3), 2001, pp. 225-239.

[25] Constant, D., Sproull, L.S., and Kiesler, S., "What's Mine Is Ours, or Is It? A Study of Attitudes About Information Sharing", *Information Systems Research*, 5(4), 1994, pp. 400-421.

[26] Carillo, P., Robinson, H., Al-Ghassani, A., and Anumba, C., "Knowledge Management in Uk Construction: Strategies, Resources and Barriers", *Project Management Journal*, 35(1), 2004, pp. 46-56.

[27] Kumaraguru, P., Sheng, S., Acquisti, A., and Cranor, L.F., "Lessons from a Real World Evaluation of Anti-Phishing Training", in (Editor, 'ed.'^eds.): Book Lessons from a Real World Evaluation of Anti-Phishing Training, IEEE, 2008, pp. 1-12.

[28] Wright, R.T., Chakraborty, S., Basoglu, A., and Marett, K., "Where Did They Go Right? Understanding the Deception in Phishing Communications", *Group Decision and Negotiation*, 19(4), 2010, pp. 391-416.

[29] Pavlou, P.A., and Gefen, D., "Building Effective Online Marketplaces with Institution-Based Trust", *Information Systems Research*, 15(1), 2004, pp. 37-59.

[30] Jarvenpaa, S.L., Tractinsky, N., and Saarinen, L., "Consumer Trust in an Internet Store: A Cross-Cultural Validation", *Journal of Computer-Mediated Communication*, 5(2), 1999,

[31] Malhotra, N.K., Kim, S.S., and Agarwal, J., "Internet Users' Information Privacy Concerns (Iuipc): The Construct, the Scale, and a Causal Model", *Information Systems Research*, 15(4), 2004, pp. 336-355.

[32] Thatcher, J.B., Zimmer, C., Gundlach, M.J., and Mcknight, D.H., "Internal and External Dimensions of Computer Self-Efficacy: An Empirical Examination", IEEE Transactions on Engineering Management, 55(4), 2008, pp. 628-644.

[33] Grover, V., Lee, C.C., and Durand, D., "Analyzing Methodological Rigor of Mis Survey Research from 1980-1989", Information & Management, 24(1993, pp. 305-317.

[34] Brown, T.A., Confirmatory Factor Analysis for Applied Research, The Guilford Press, New York, NY, 2006.

[35] Brown, T.A., Confirmatory Factory Analysis for Applied Research, The Guilford Press, New York, NY, 2006.

[36] Fornell, C., and Larcker, D.F., "Evaluating Structural Equations Models with Unobservable Variables and Measurement Error", Journal of Marketing Research, 18(1), 1981, pp. 39-50.

[37] Wasko, M.L., and Faraj, S., "Why Should I Share? Examining Social Capital and Knowledge Contribution in Electronic Networks of Practice", MIS Quarterly, 29(1), 2005, pp. 35-57.

[38] Constant, D., Sproull, L., and Kiesler, S., "The Kindness of Strangers: The Usefulness of Electronic Weak Ties for Technical Advice", Organization Science, 7(2), 1996, pp. 119-135.

[39] Bock, G.-W., Zmud, R.W., Kim, Y.-G., and Lee, J.-N., "Behavioral Intention Formation in Knowledge Sharing: Examining the Roles of Extrinsic Motivators, Social-Psychological Forces, and Organizational Climate", MIS Quarterly, 2005, pp. 87-111.

[40] Burt, R.S., Structural Holes: The Social Structure of Competition, Harvard university press, 2009.

[41] Meservy, T.O., Jensen, M.L., and Fadel, K.J., "Evaluation of Competing Candidate Solutions in Electronic Networks of Practice", Information Systems Research, 25(1), 2014, pp. 15-34.

[42] Herath, T., and Rao, H.R., "Protection Motivation and Deterrence: A Framework for Security Policy Compliance in Organisations", European Journal of Information Systems, 18(2), 2009, pp. 106-125.

[43] Johnston, A.C., and Warkentin, M., "Fear Appeals and Information Security Behaviors: An

Empirical Study", MIS Quarterly, 34(3), 2010, pp. 549-566.

Appendix A - Measurement Model

	AVE	C.R.	1	2
1 - RISK	0.53	0.82	0.73	
2 - TRUST	0.70	0.91	0.14	0.84
3 - PHEXP	0.65	0.88	0.04	-0.03
4 - CSE_INT	0.61	0.82	0.06	0.15
5 - CSE_EX	0.64	0.84	-0.10	0.18
6 - MOTIV	0.77	0.91	-0.04	-0.03
	3	4	5	6
1 - RISK				
2 - TRUST				
3 - PHEXP	0.81			
4 - CSE_INT	0.39	0.78		
5 - CSE_EX	0.20	0.32	0.80	
6 - MOTIV	0.08	0.02	0.23	0.88

Construct	Items	Loadings
RISK	RISK1	0.759
	RISK2	0.681
	RISK3	0.766
	RISK4	0.710
TRUST	TRUST1	0.860
	TRUST2	0.823
	TRUST3	0.846
	TRUST4	0.825
PHEXP	IDPHISH1	0.872
	IDPHISH2	0.967
	IDPHISH3	0.655
	IDPHISH4	0.687
CSE_INT	CSE_INT1	0.854
	CSE_INT2	0.787
	CSE_INT3	0.689
CSE_EXT	CSE_EX1	0.675
	CSE_EX2	0.982
	CSE_EX3	0.698
MOTIVATION	MOTIV_1	0.779
	MOTIV_2	0.962
	MOTIV_3	0.875