

---

Research paper

# Categorizing human phishing difficulty: a Phish Scale

Michelle Steves, Kristen Greene\* and Mary Theofanos

National Institute of Standards and Technology Gaithersberg, MD 20899, USA

\*Corresponding address, Kristen Greene, National Institute of Standards and Technology, 100 Bureau Drive, Gaithersberg, MD 20899, USA. Tel: 301-975-8119; E-mail: kristen.greene@nist.gov

Received 2 August 2019; revised 4 December 2020; accepted 17 April 2020

## Abstract

As organizations continue to invest in phishing awareness training programs, many chief information security officers (CISOs) are concerned when their training exercise click rates are high or variable, as they must justify training budgets to organization officials who question the efficacy of awareness training when click rates are not declining. We argue that click rates should be expected to vary based on the difficulty of the phishing email for a target audience. Past research has shown that when the premise of a phishing email aligns with a user's work context, it is much more challenging for users to detect a phish. Given this, we propose a Phish Scale, so CISOs and phishing training implementers can easily rate the difficulty of their phishing exercises and help explain associated click rates. We base our scale on past research in phishing cues and user context, and apply the scale to previously published and new data from enterprise-based phishing exercises. The Phish Scale performed well with the current phishing dataset, but future work is needed to validate it with a larger variety of phishing emails. The Phish Scale shows great promise as a tool to help frame data sharing on phishing exercise click rates across sectors.

**Key words:** phishing defences; embedded phishing awareness training; Phish Scale; cybersecurity defences; phishing cues; phishing email premise

---

## Introduction

According to Cybersecurity Ventures' 2019 Official Annual Cybercrime Report, it is estimated that cybercrime damages will cost the world \$6 trillion annually by 2021 [1]. These cost projections are supported by historical cybercrime figures and recent year-over-year growth. Furthermore, the report relates there has been a notable increase in hacking activities sponsored by hostile nation states, as well as activities from organized crime syndicates. Finally, the cyber attack surface, the sum of vulnerabilities that can be exploited to carry out an attack, continues to grow, in large part due to an explosion of Internet of Things (IoT) devices. Humans are another particularly important component of the overall attack surface, as social engineering continues to be successful. In recognition of the importance of human behavior in cybersecurity, organizations are more widely investing in cybersecurity awareness programs for

their computer users, with significant focus on phishing training. Embedded phishing awareness training is popular—and in some cases, mandated—in a wide variety of sectors, including financial services, government, healthcare, and academia. In this type of training, simulated phishing emails are sent that mimic real-world threats to raise employee phishing awareness.

Many Chief Information Security Officers<sup>1</sup> (CISOs) have expressed concern when their training exercise click rates are high. This is especially true for more mature or long-running awareness programs, as CISOs often expect progressively lower click rates to show the effectiveness of training. Further, the Return on Investment (ROI) for such training may be questioned if click rates are high or even variable. However, low click rates do not necessarily indicate training effectiveness and may instead mean the phishing emails used were: (i) too easy, (ii) not contextually relevant for most

---

<sup>1</sup> Also referred to as "Senior Agency Information Security Officer" (SAISO), NIST SP 800-37, Rev. 2 [45].

staff, or (iii) the phish was repeated or very similar to previous exercises. In fact, low click rates and training programs in general, can generate a false sense of security or complacency if considered in isolation. Phishing awareness training program click rates must be part of a more comprehensive, metrics-informed approach to effectively understand and combat phishing threats [2].

Past work [3] has shown that click rates will vary based on the contextual relevance of the phish, with highly contextually relevant phish resulting in extreme spikes in click rates—despite years of phishing awareness training. Furthermore, attackers continue to refine and vary phishing attack premises. Although “traditional” phishing emails that focus on credential harvesting are still quite successful, attackers are becoming more sophisticated and creative all the time, this includes having added malware delivery. Additionally, there is a treasure trove of readily available information online that attackers can use to better tailor phish and capitalize on contextual relevance. While some information is willingly and openly shared by users on social media, much other information has been exposed through large-scale data breaches, such as the recent Facebook hack [4].

While repetition is important for training phishing recognition and for conditioning reporting behavior, simple repetition of the same or very similar phishing emails does not represent the full spectrum of phishing threats observed in the real world. It is important to vary phishing exercises appropriately and challenge staff with contextually relevant phish of varying difficulty to provide training on new scams for which variable click rates should be expected. This should not be viewed as a negative effect but rather a positive outcome, as it means organizations are truly training their staff with phish that represent current real-world threats. But how exactly does one measure the difficulty of a given phishing email? While we can certainly measure click rates post-hoc and infer detection difficulty somewhat from those numbers, we would prefer an a priori method of difficulty determination. In discussions with CISOs at Healthcare CyberGard Annual Conference [5], Information Security and Privacy Board [6], and others, we found that a method to determine phishing message difficulty would indeed be highly beneficial for those responsible for phishing training implementation. To meet this need, we propose a Phish Scale, an easy way for CISOs and training implementers to characterize the difficulty of their phishing exercises and provide context for the associated metrics. This context is a missing element that training implementers need to improve the training benefit of their exercises and subsequent ROI.

In this article, we describe our exploratory effort to construct a preliminary conceptual version of the Phish Scale and its components. Further, we use the Phish Scale to determine the difficulty rating for each of ten real-world phishing training exercises that are described in detail. Then we observe if the Phish Scale difficulty rating for each exercise aligns with the exercise’s actual click rate. Finally, we discuss our observations, limitations of the effort to date, as well as future work including refinement of its components and validation of the overall scale through a wider variety of phishing emails.

## Background

Technological and human-centered approaches are used in conjunction to combat email phishing. Technologically focused approaches include mechanisms like filtering, firewalls, and blacklists, whereas human-centered approaches tend to focus on cybersecurity awareness training, and often on phishing specifically. Due in part to

advances in Machine Learning (ML) and Artificial Intelligence (AI), email filters in particular, are becoming ever more effective at blocking generic spam. This has meant that users now see fewer emails of this nature in their inboxes. Recent work has posited the existence of the “Prevalence Paradox” [7], suggesting that users may therefore be more vulnerable when such emails do get through, due to their reduced experience with potentially malicious emails. Yet other work by Levari *et al.* [8] has shown that people often expand their concept of a given stimulus in response to a decrease in the prevalence of said stimulus, for example, seeing neutral faces as threatening when threatening faces became rare. Although the series of experiments by Levari *et al.* did not address phishing specifically, given the set of topics they investigated, it would certainly be plausible to expect their findings to hold in the phishing domain. We hope additional research on the effects of prevalence on phishing detection—for both humans and AI—will reconcile different findings on prevalence.

In addition to prevalence, there are numerous other factors that complicate human detection of phishing emails. There are several existing theories and models of phishing susceptibility that are highly relevant for the development of a Phish Scale. These theories and models directly address the types of email cues, tactics, and individual user characteristics that together help—at least partially—explain the relative ease or difficulty of human phishing detection.

Protection Motivation Theory or PMT [9] addressed user perceptions of threat and corresponding perceived threat management ability. PMT has largely been applied to security behavior in general, although Wang *et al.* [10] did apply PMT specifically to phishing threat perception. Much more recently than PMT, which was originally proposed in 1975, an Integrated Information Processing Model of Phishing Susceptibility, or IIPM, was proposed [11]. The IIPM proposed that users’ limited attentional resources for information processing are essentially hijacked when certain techniques like urgency are used to influence behavior, meaning that users rely on heuristic information processing (System I [12]), rather than engaging in deeper, more systematic processing (System II [12]). When this type of surface-level information processing style is used it makes users more likely to overlook or ignore cues that might otherwise tip a user off as to the legitimacy of the email, such as an incorrect sender address. In 2016, Vishwanath *et al.* proposed the Suspicion, Cognition, and Automaticity Model, or SCAM, which posited that individual user characteristics cause variability in the use of heuristic processes for email evaluation [13].

Recent work by Williams, Hinds, and Joinson [14] considered these three models (PMT, IIPM, and SCAM) within the work context of an international organization with sites in the UK, finding that the presence of authority cues increased the likelihood that users would click a suspicious email link. In addition to the types of models or theories such as PMT, IIPM, SCAM, there is a large wealth of prior work investigating or describing the impact of particular email cues, such as inclusion of authority and urgency cues. Research on phishing cues is particularly relevant for development of a Phish Scale, as email users rely on cues to determine if a particular email message is a phish.

Indeed, anti-phishing advice and training stress the characteristics of phishing messages that email users should look for; these are often called cues, indicators and hooks. The list of cues is long and varied, such as those contained in refs [15, 16]. Because there is no set pattern of which cues may be contained in any particular message, the task for users when determining if a message is a phish is harder than if the list were very short. Making the task even more difficult, prior work shows the alignment of the phish’s premise and

user context affects which cues the user finds to be salient. Further, the same cue can be compelling for some users but suspicion generating for others—depending on the user’s context [2, 3].

In the Greene *et al.* [3] study, phishing exercise data were collected over 4.5 years in an ecologically valid workplace setting, with corresponding survey data for the final year. The study found that user context was extremely important in phishing susceptibility; the authors proposed that it was the lens through which users viewed and interpreted email cues. When a user’s work context was misaligned with the premise of the phishing email, they were more likely to attend to suspicious cues, for example, they have no invoicing responsibilities at work and the phishing email was purportedly an unpaid invoice. In contrast, when a user’s work context was well aligned with the phishing email premise, they were more likely to attend to compelling cues, and completely ignore or largely discount suspicious cues. In this case, if the user is directly responsible for paying invoices at work and the phishing email was purportedly an unpaid invoice.

Greene *et al.* [3] emphasized the importance of phishing research in the workplace setting, as much prior phishing work was conducted in laboratories with artificial user contexts or university settings that can be quite different than the workplace. Williams *et al.* [14] also addressed this need for workplace data in their research. One of the few other studies situated in the workplace was conducted by Caputo *et al.* [17], but due to limitations was only able to suggest the possible importance of user context. We further contribute to the growing corpus of workplace-based phishing research, by applying our Phish Scale to three previously published workplace-based phishing exercises in ref. [3], four phishing exercises detailed in ref. [18], as well as, reporting on and applying the Phish Scale to data from three previously unpublished workplace-situated phishing exercises. Eight of the ten exercises have *n*’s of  $\sim 70$  for each exercise, while two of our exercises have much larger sample sizes, with *n*’s of  $\sim 5\,000$  for each exercise.

## Method

To assist organizations tasked with implementing phishing awareness training programs, it is important to consider the relative detection difficulty of training messages. Phishing messages, whether those intended for training or actual threats, can be more or less difficult for a given work group to detect as a phishing attempt. Understanding the detection difficulty helps phishing awareness training implementors in two primary ways: (i) by providing context regarding training message click and reporting rates for a target audience, and (ii) by providing a way to characterize actual phishing threats so the training implementor can reduce the organization’s security risk by tailoring training to the types of threats their organization is facing. To this end, we developed the Phish Scale to help practitioners rate detection difficulty of both training and actual threat phishing messages. The Phish Scale is intended to contextualize click rates for embedded phishing awareness training as well as tailor training efforts. We anticipate it will provide CISOs with another metric to help gauge the progress of their awareness programs over time and address risk. The scale is intended to categorize the detection difficulty of a phishing message with respect to a target audience.

In this section, we describe the Phish Scale and the operationalization of its components into a single framework. In the next section, we present data from ten workplace-situated phishing

awareness training exercises to illustrate how to derive a phish difficulty rating using the Phish Scale.

## The Phish Scale

To develop our Phish Scale, we began by considering the primary elements that CISOs and training implementors use when selecting and customizing phishing training exercises. These elements are scenario premise and message content. The scenario premise may pertain to a relatively new threat or an older threat that remains effective for a particular target audience. The message content is typically customizable by the trainer and contains the cues that trainees might use to detect the training phish. For this exploratory effort, we root the Phish Scale in these two primary elements: the *cues contained in the message* and the *premise alignment for the target audience*.

Other factors such as personality, curiosity, distractedness, concern for security, and the like certainly affect click rates, and ultimately we intend to consider how to account for additional factors such as these; we return to this topic in the future work section. However, for now, this effort starts with message cues and premise alignment as these elements undoubtedly play crucial roles in phishing detection by humans and, importantly, they can be categorized by training implementors for a given target audience. For this initial effort at characterizing detection difficulty, the Phish Scale components are:

1. A rating system for observable characteristics of the phishing email itself, such as the number of cues, nature of the cues, repetition of cues, and so on.
2. A rating system for alignment of the phishing email premise with respect to a target audience.

Table 1 presents our exploratory, conceptual framework illustrating how detection difficulty rating is derived once the categories for number of cues and premise alignment are determined. In an attempt to keep the categorization relatively simple for training implementors, we used three categories for each component and assigned labels representing relative ranges for each. Briefly, the three categories used to reflect the effect of the “cues” contained in a phishing message are *few* (lower count equating to fewer opportunities to detect), *some* (moderate number of cues), and *many* (higher number of cues, more opportunities to detect). Similarly, three categories are used to characterize the premise alignment; they are *high*, *medium*, and *low*. We discuss the scale components in more depth in the next sections after a few observations about the conceptual framework.

In the conceptual framework we acknowledge the stronger influence of premise alignment component over cues; this is consistent with findings reported in ref. [3]. The stronger premise alignment influence is reflected in the detection difficulty rating tending to be at

**Table 1:** the Phish Scale

Number of cues	Premise alignment	Detection difficulty
Few (more difficult)	High	Very difficult
	Medium	Very difficult
	Low	Moderately difficult
Some	High	Very difficult
	Medium	Moderately difficult
	Low	Moderately to Least difficult
Many (less difficult)	High	Moderately difficult
	Medium	Moderately difficult
	Low	Least difficult

the *Very difficult* or *Moderately difficulty* rating when the premise alignment is categorized as *High* or *Medium*. Additionally, there are more *Very difficult* detection difficulty rating assignments than *Least difficult* rating assignments in the entire conceptual Phish Scale framework.

The detection difficulty rating for the combination of *Some* cues and *Low* premise alignment was given a range from *Moderately to Least difficult* rating, further reflecting our belief that even a *Low* premise alignment can have a disproportionate effect on increasing detection difficulty. While we expect all of the ratings to be informed with empirical data, this is especially true for this particular combination (low premise alignment and some cues). Finally, we purposefully did not label a category as *Easy to detect* or similar, as we expect that the premise of any phishing message will typically align for *at least a few users* and for them, detection is often not easy.

Ultimately, we anticipate that each detection difficulty rating will equate to a range of click rates. For example, the phishing training messages that have a corresponding detection difficulty rating of *least difficult* may be expected to have a click rate of less than 10%. We return to this topic in the discussion.

Next we discuss the operationalization of each scale component. In the following section, we walk through marrying the real-world phishing training exercise data with these conceptual categorizations and discuss our observations.

#### Phishing message cues

To incorporate the effect of phishing message cues in the scale, we decided to use the count of instances of those characteristics that are present in the message being rated. Our reasoning is that the fewer phishing cues present in a message, the more difficult it is to detect. Conversely, the more cues present, the more opportunities for a user to notice a tip-off that generates suspicion. We realize the effect of any single cue or hook can differ from instance to instance and person to person. Indeed, we return to this topic in the discussion. Currently, there are three categories in the framework to describe the quantity of these characteristics: *Few* (fewer opportunities to detect), *Some*, and *Many* (more opportunities to detect).

Before we can count cues, we needed to determine which phishing characteristics—the list of cues, indicators and hooks—are appropriate for inclusion in the framework. From ref. [3], we see that a particular phishing characteristic may either be suspicion-generating (a tip-off) or compelling (a hook), depending on the user's context. In keeping with prior literature, we use the term "cue." However, we mean it in the broader sense of a phishing message characteristic. We require that each cue included in the framework be able to be tied to an objectively observable characteristic in a message.

From the literature, we considered the compendiums of phishing cues in refs [15] and [16]. We used the cues given in ref. [15] as a starting point. Additionally, we modified the categories in an attempt to order the cues from those that are often suspicion-generating, such as errors, to those that are typically compelling, such as common tactics, these tactics being commonly used because they continue to be compelling. This is a rough ordering of categories at best, but we felt it is better suited to counting cues than those given in refs [15] and [16]. The categories are: *Error*—relating to spelling and grammar errors and inconsistencies contained in the message; *Technical indicator*—pertaining to email addresses, hyperlinks and attachments; *Visual presentation indicator*—relating to branding, logos, design and formatting; *Language and content*—

such as a generic greeting and lack of signer details, use of time pressure and threatening language; and, *Common tactic*—use of humanitarian appeals, too good to be true offers, time-limited offers, poses as a friend, colleague, or authority figure, and so on. Beyond the cues and characteristics given in [15], we wove in additional phishing characteristics from refs [3, 16] and many others.

Table 2 provides the list of cues we identified that are objectively present in phishing messages. Further, it also contains a brief description of each, associated references, and the criteria we used when deciding if a particular cue was observably present in an individual message. To determine the cues count, use the criteria for each cue, count how many instances for each and sum for a total.

For this initial effort, we recognize this list is not exhaustive and will be expanded. Additionally, we anticipate some form of weighting may be useful to reflect cue saliency. However, given the variability in cue saliency for individuals within a target population, this is a non-trivial exercise. We expect such refinements will be closely examined with additional development of the scale.

For the purpose of the Phish Scale, we did not include phishing message cues related to mismatches with the user's world, such as an individual's particular work responsibilities or an individual's expectations, for example expecting an important phone call. Work responsibilities and general workplace expectations for the target audience are folded into the premise alignment component of the Phish Scale, described next.

#### Phishing premise alignment

Incorporating premise alignment is a process of characterizing the pertinence of the email message premise for the target audience. It attempts to capture alignment with the following for a target audience: work responsibilities and business practice plausibility, workplace pertinence of the topic, alignment with other events or situations, including external to the workplace, concern over not clicking, and exposure or warnings about the premise that would affect the tendency to click. Another way to view alignment is what makes the premise compelling. Overall, we are attempting to categorize premise alignment, not premise misalignment, which in some ways are the reverse of each other. Even so, we acknowledge the mitigating effect of training and awareness on phishing recognition and believe it should be considered when categorizing premise alignment. A particularly good example of this mitigating effect is the wide-spread awareness of the Nigerian 419 scam [24]—which is so well-known that its many varieties are typically recognized as scams.

We expect premise alignment will be determined by the training implementor for a particular phishing message—*someone with knowledge of the target audience's work culture, responsibilities and expectations as a group*. Premise alignment cannot be determined in the abstract; knowledge of the target population's context of work with respect to the phishing message's premise is crucial in accurately categorizing premise alignment. We use three categories to characterize the alignment: *High*, *Medium*, and *Low*. To determine premise alignment, the training implementor must understand and categorize the premise alignment for the target audience. We developed two methods to categorize premise alignment. One uses a blended perspective, while the second uses a formulaic approach to premise alignment categorization. The blended perspective method uses three ratings, high, medium, and low, to categorize how strongly the premise aligns for portions of the target audience in broad terms and then melds those into an overall rating. The formulaic approach has five elements that are each rated and then the ratings are

**Table 2:** operationalized phishing message cues

Cue type	Cue name	Description	References	Criteria for counting	
Error	Spelling and grammar irregularities	Spelling or grammar errors, mismatched plurality and so on	[11, 19–30]	Does the message contain spelling or grammar errors, including mismatched plurality?	
	Inconsistency	Inconsistent content within the email	[3]	Are there inconsistencies contained in the email message?	
Technical indicator	Attachment type	The presence of file attachments, especially an executable	[31]	Is there a potentially dangerous attachment?	
	Sender display name and email address	Spoofed display names - hides the sender and reply-to email addresses	[11, 13, 19, 21, 22, 24, 27, 29, 32]	Does a display name hide the real sender?	
	URL hyperlinking	URL hyperlinking hides the true URL behind text; the text can also look like another link	[20–22, 25, 27, 33]	Is there text that hides the true URL in a hyperlink?	
	Domain spoofing	Domain name used in email address and links looks similar to plausible	[3, 34]	Is a domain name used in addresses or links plausibly similar to a legitimate entity's domain?	
Visual presentation indicator	No/minimal branding and logos	No or minimal branding and logos	[13, 19, 22, 23, 25, 27, 32, 34, 42]	Is appropriate branding missing?	
	Logo imitation or out-of-date branding/logos	Spoof or imitation of logo/out-of-date logo	[3, 24]	Do any branding elements appear to be an imitation or out-of-date?	
	Unprofessional looking design or formatting	Formatting and design elements that do not appear to have been professionally generated	[25, 27, 28, 34–36]	Does the design and formatting violate any conventional professional practices?	
Language and content	Security indicators and icons	Security indicators and icons	[25, 35]	Are any inappropriate security indicators or icons present?	
	Legal language/copyright info/disclaimers	Any legal type language such as copyright information, disclaimers, tax implications	[25]	Does the message contain any legal type language such as copyright information, disclaimers, tax information?	
	Distracting detail	Distracting Detail	[3]	Does the message contain any detailed aspects that are not central to the content?	
	Requests for sensitive information	Requests for sensitive information, like a Social Security number or other identifying information	[3, 21, 22]	Does the message contain a request for any sensitive information, including personally identifying information or credentials?	
	Sense of urgency	Use of time pressure to try to get users to quickly comply with the request	[11, 20–22, 24, 27, 32, 37]	Does the message contain time pressure, including implied?	
	Threatening language	Use of threats such as legal ramifications	[11, 20, 21, 27, 32, 37]	Does the message contain a threat, including an implied threat?	
	Generic greeting	A generic greeting and an overall lack of personalization in the email	[20, 21, 24, 27, 28, 33, 34, 37]	Does the message lack a greeting or lack personalization in the message?	
	Lack of signer details	Emails including few details about the sender, such as contact information	[24, 32]	Does the message lack detail about the sender, such as contact information?	
	Common tactic	Humanitarian appeals	Appeals to help others in need	[24, 27, 32]	Does the message make an appeal to help others?
		Too good to be true offers	Contest winnings or other unlikely monetary and/or material offerings	[23, 24, 27, 28, 30]	Does the message offer anything that is too good to be true, such as having won a contest, lottery, free vacation and so on?
You're special		Just for you offering... such as a valentine e-card from a secret admirer	[24]	Does the message offer anything just for you?	
Limited time offer		This offer won't last long...	[24]	Does the message offer anything for a limited time?	
Mimics a work or business process – a legitimate email		Mimics any plausible work process such as new voicemail, package	[24]	Does the message appear to be a work or business-related process?	

**Table 2.** (continued)

Cue type	Cue name	Description	References	Criteria for counting
		delivery, order confirmation, notice of invoice, and so on		
	Poses as friend, colleague, supervisor, authority figure	Email purporting to be from a friend, colleague, boss or other authority figure	[14, 24]	Does the message appear to be from a friend, colleague, boss or other authority entity?

summed to obtain an overall rating. One or the other method may appeal more to individual training implementors; we describe both.

*Method 1: Blended perspective premise alignment categorization.* A rating is chosen based on the following guidelines.

1) **High alignment.** For high premise alignment, there should be a significant portion of the target audience for which the premise matches work responsibilities or practices, is highly plausible, and/or aligns strongly with an audience-relevant event. For example, if the recipient population is the finance department and the phishing message has a premise of a late or missed payment, the overall alignment is high.

2) **Medium alignment.** Medium alignment is achieved with either case: (i) when the premise has plausible but weak context alignment with a large portion of the target audience or (ii) when the premise has moderate context alignment with a small portion of the target audience. For example, if the recipient population mostly works in one physical location and the phishing message has a moderately pertinent premise for the few members of the recipient population who work in another physical location.

3) **Low alignment.** There is low alignment when the premise pertains to a topic that is not relevant or plausible to the target audience. For example, if the recipient population is the finance department and the phishing message premise pertains to a Call for Papers on biotech research or a similarly unrelated topic, the overall alignment is low.

*Method 2: Formulaic approach to premise alignment categorization.* This approach uses a set of five elements that are each rated on a 5-point scale. The overall score determined from the ratings yields a value reflecting a premise alignment categorization. These elements were chosen because they are aspects relating to premise alignment that the training implementer can categorize for the target training population. The message premise elements are:

1. *Mimics a workplace process or practice:* this element attempts to capture premise alignment with workplace process or practice for the target audience,
2. *Has workplace relevance:* this element attempts to reflect pertinence of the premise for the target audience,
3. *Aligns with other situations or events, including external to the workplace:* alignment with other situations or events, even those external to the workplace lends an air of familiarity to the message,
4. *Engenders concern over consequences for NOT clicking:* potentially harmful ramifications for not clicking raise the likelihood to clicking [3],

5. *Has been the subject of targeted training, specific warnings, or other exposure:* this element is intended to reflect targeted training effects that would lead to premise detection. Care must be taken to appropriately incorporate the training or warning specificity, as transfer of learning is quite difficult [38].

We use the following 5-point rating scale of even numeric values of zero to eight with these associated anchors:

- 8 = **Extreme** applicability, alignment, or relevancy
- 6 = **Significant** applicability, alignment, or relevancy
- 4 = **Moderate** applicability, alignment, or relevancy
- 2 = **Low** applicability, alignment, or relevancy
- 0 = **Not applicable**, no alignment, or no relevancy

The overall premise alignment score for a particular phish and its target audience is the sum of the ratings of elements 1 through 4. Since the fifth element pertains to training on the premise and helps with detection, its score is subtracted from overall sum. The highest score possible is 32, indicating very high premise alignment. The lowest score possible is -8, showing extremely low premise alignment.

## Application of the Phish Scale

In this section, we present data from ten phishing training exercises and use the Phish Scale to determine the detection difficulty rating for each exercise. The exercise data used here originated from a project that was started in 2012 by the Information Technology Security and Networking Division (ITSND) at NIST as a long-term trial deployment of an embedded awareness training effort. The trial deployment was intended as a multi-year effort and used a commercially available system to help develop and deliver phish messages and training, as well as track click rates. The same system was used throughout the entire 5+ year period.

For all but the last two exercises conducted in the trial, the targeted population within the institute was one operating unit (OU) having approximately 70 staff members. The awareness training provided by these exercises augmented the IT security awareness training the entire institute's workforce received annually. OU staff were aware their unit was participating in the trial. These exercises were conducted by the OU's Information Technology Security Officer (ITSO). The same person held the position during the entire trial period. The ITSO selected the phishing message and its premise from templates provided by the training system that mimicked current real-world threats. Some messages were tailored to align with business and communication practices within the organization or were personalized, in other words, they were spear phish [39]. For the two remaining training exercises, the entire NIST staff was the target population. For all exercises, the phishing training emails modeled real-world phishing campaigns, and participating staff were in their normal work environments with their regular work

**Table 3:** safety requirements exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	8 (Extreme)
Has workplace relevance	8 (Extreme)
Aligns with other situations or events, including external to the workplace	8 (Extreme)
Engenders concern over consequences for NOT clicking	8 (Extreme)
Has been the subject of targeted training, specific warnings, or other exposure	-2 (Low)
Overall alignment	30

loads, providing ecological validity. All exercises were unannounced and deployed at irregular intervals to avoid priming effects. The data were gathered with appropriate human subjects approval at the National Institute of Standards and Technology (NIST).

First, we provide a description of each exercise, its premise alignment rationale (Method 1), premise alignment element ratings and overall score (Method 2), and a brief description of the target audience and size. Then we gather the data in Table 13, including the observed cue counts provided in Supplementary Appendix A, and show the detection difficulty rating for each exercise side-by-side with the actual click rate.

Each exercise has been given a label pertaining to the scenario, such as Safety requirements, as well as an exercise number, E1, E2, . . . , and so on, pertaining to the presentation order, which is ordered by click rate. Note that although we assigned premise alignment category ratings to these exercises—rather than the training implementers on their own—we did so with pertinent input from the training implementers. An image of each phishing message is provided in Supplementary Appendix B. These images are recreations of the actual message screen captures, although the message recipients have been replaced with fictitious names.

Exercises E3, E5, and E7—New voicemail, Unpaid invoice, and Order confirmation, respectively, were initially reported in Greene *et al.* [3]. Exercises E1, E2, E6, and E8 – Safety requirements (formerly labeled ‘Gmail’), Weblogs, Valentine, and Security token, respectively, were originally presented in ref. [18]. Exercises E4, E9, and E10—Scanned file, Gift certificate, and Adobe update, respectively, represent new data.

### Phishing exercise descriptions

#### Safety requirements<sup>2</sup> (E1)

*Message description:* The safety requirements phish was a particularly clever spear phish. It targeted employees using a spoofed upper management Gmail address, a tactic based on a real-world phish previously observed at NIST. The real-world phish appeared to come from the personal Gmail account of NIST’s director (firstname.lastname1@gmail.com) and went to a list of laboratory managers. The training exercise phish also appeared to come from NIST upper management, the organizational unit (OU) director, and appeared as (firstname.lastname1@gmail.com). The similarities continued in the subject line, body, and closing portions of the email. The subject line was, “PLEASE READ THIS,” which is important given that NIST has a very strong emphasis on fostering a culture of safety. The email was personalized with the recipient’s first name.

2 In [18], we had labeled this exercise ‘Gmail.’ We changed the label to represent the premise rather than the tactic to be consistent with other exercise labels we use.

**Table 4:** weblogs exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	8 (Extreme)
Has workplace relevance	8 (Extreme)
Aligns with other situations or events, including external to the workplace	0 (Not applicable)
Engenders concern over consequences for NOT clicking	8 (Extreme)
Has been the subject of targeted training, specific warnings, or other exposure	0 (Not applicable)
Total	24

The body said, “I highly encourage you to read this.” The next line contained a link with the following bolded text, “Safety Requirements.” The email was signed simply with, “Best regards,” and the first name of the OU’s director.

*Premise alignment (Method 1):* The alignment is categorized as High—the premise alignment is very strong given the larger organization’s substantial emphasis on workplace safety and that the message appeared to come from upper NIST management—notable as an authority figure in this context. Alignment is further strengthened by the target department’s responsibility for the larger organization’s occupational health and safety.

*Premise alignment (Method 2):* The overall alignment is 30 out of a possible 32. Table 3 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One OU within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support,  $n = 73$ .

#### Weblogs (E2)

*Message description:* The weblogs phish was another spear phish. It appeared to come from a system administrator with the email address, notice@nist.gov. The subject line was, “Unauthorized Web Site Access.” There was no personalization. The body said, “\*This is an automated email\* Our regulators require we monitor and restrict certain website access due to content. The filter system flagged your computer as one that has viewed or logged into websites hosting restricted content. The system is not fool-proof and may incorrectly flag restricted content. The IT department does not investigate every web filter report, but disciplinary action may be taken.” In bold, it said, “Log into the filter system with your network credentials immediately and review your logs to see which websites triggered this alert.” This was followed by a link that was labeled, “Web Security Logs.” There was no contact information given, and the email closed with, “Do not reply to this email. This email was automatically generated to inform you of a violation of our security and content policies.”

*Premise alignment (Method 1):* The alignment is categorized as High—the premise aligns with the fact that accessing inappropriate content is indeed a violation of the organization’s Rules of Conduct policy and can be grounds for dismissal for anyone at the organization. The premise capitalizes on the fact that many organizations, including NIST, scan log data routinely. The threat component coupled with the severity of the consequences increases the alignment. Of note, all new employees receive in-person training regarding the organization’s Rules of Conduct and Information

**Table 5:** unpaid invoice exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	4 (Moderate)
Has workplace relevance	8 (Extreme)
Aligns with other situations or events, including external to the workplace	6 (Significant)
Engenders concern over consequences for NOT clicking	6 (Significant)
Has been the subject of targeted training, specific warnings, or other exposure	0 (Not applicable)
Total	24

**Table 6:** scanned file exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	6 (Significant)
Has workplace relevance	4 (Moderate)
Aligns with other situations or events, including external to the workplace	6 (Significant)
Engenders concern over consequences for NOT clicking	4 (Moderate)
Has been the subject of targeted training, specific warnings, or other exposure	-2 (Low)
Total	18

Technology (IT) policies, where the disciplinary actions associated with inappropriate web content viewing are highly stressed.

*Premise alignment (Method 2):* The overall alignment is 16 out of a possible 32. Table 9 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One OU within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support,  $n = 64$ .

#### Unpaid invoice (E3)

*Message description:* The unpaid invoice phish appeared to be from a fictitious employee of the same institution as the email recipients, a fellow Federal employee named Jill Preston (jill.preston@nist.gov). The subject line was, "Unpaid invoice #4806." The greeting was personalized with "Dear [Firstname Lastname]." The email body said, "Please see the attached invoice (.doc) and remit payment according to the terms listed at the bottom of the invoice. Let us know if you have any questions. We greatly appreciate your prompt attention to this matter!" The email simply closed with the name "Jill Preston." There was no other contact information included below the name. Of note, there was a file extension mismatch between the way the attachment was referred to in the body of the email (as a .doc) and the way the attachment itself was labeled, it appeared to be a .zip, with the filename, "invoice\_S-37644806.zip". The unpaid invoice phish mimicked the Locky ransomware [39], a real-world threat current at that time.

*Premise alignment (Method 1):* The alignment is categorized as High—the premise aligned extremely highly for roughly a third of the target audience and aligned somewhat for the remainder of the department. Additionally, the whole of the targeted OU was on alert

**Table 7:** new voicemail exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	6 (Significant)
Has workplace relevance	4 (Moderate)
Aligns with other situations or events, including external to the workplace	2 (Low)
Engenders concern over consequences for NOT clicking	4 (Moderate)
Has been the subject of targeted training, specific warnings, or other exposure	-2 (Low)
Total	14

for any unpaid invoices following a recent event surrounding a legitimate unpaid invoice.

*Premise alignment (Method 2):* The overall alignment is 24 out of a possible 32. Table 5 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One OU within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support.  $n = 73$

#### Scanned file (E4)

*Message description:* The scanned file phish appeared to come from "LaserPro\_2\_2\_e" with the email address, "laerpro\_2\_2\_e@nagts.org." The subject line was, "Scan from Laser Pro i780 Second Floor." There was no personalization. The body contained the text, "Please open the attachment. It was scanned and sent to using a Laser Pro i789." This text was followed by, "SENT BY: INELL," "PAGES: 1," and "FILETYPE: .DOC," each on successive lines. At the bottom of the email was an image indicating an attached file and the filename, "SCN001375.doc," concluding the phish.

*Premise alignment (Method 1):* The alignment is categorized as High—the message capitalized on a common business practice using a shared scanning & printing device. Further, the OU recently had a new large, shared printer device installed. Additionally, the attached document was purportedly scanned in-house and likely seemed trustworthy, while engendering curiosity about the content and possibly about getting the scanned file to its intended recipient.

*Premise alignment (Method 2):* The overall alignment is 18 out of a possible 32. Table 6 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One OU within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support,  $n = 62$ .

#### New voicemail (E5)

*Message description:* The new voicemail phish appeared to be from a fictitious CorpVM (corpvm@webaccess-alert.com). It appeared to be a system-generated email, with the subject line reading, "You have a new voicemail." There was a large black and green banner at the top of the email with the text, "CorpVM" in white. There were no logos present in the email, however, there was a small black footer with "© 2015 CorpVM Inc." in white. The body of the email began, "You have a new voicemail!" centered in bold text, followed by, "From: Unknown Caller, Received: 03/06/2016, Length: 00:52." Below that text was a personalized [Firstname Lastname] line, followed by, "You are receiving this message because we were unable to deliver it voice message did not go through because the



voicemail was unavailable at that moment. To listen to this message, please click [here](#). You must have speakers enabled to listen to the message. \* The reference number for this message is qvfl\_cjl09-9107319601-2125579909-62. The length of transmission was approximately 52 seconds. The receiving machine’s ID: YJH35-TW410-F37JZL. Thank you.” Finally, the email closed with smaller text in italic that read, “This is a system-generated message from a send-only address. Please do not reply to this email.”

*Premise alignment (Method 1):* The alignment is categorized as Medium—the premise was plausible; around the same time as the exercise, a new business process for voicemail notification, *not* delivery, was being rolled out, although without much fanfare. Even though the premise was plausible, it had no or weak context alignment for most, although not all, of the target audience based on survey feedback reported in ref. [3].

*Premise alignment (Method 2):* The overall alignment is 16 out of a possible 32. Table 7 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One Operational Unit (OU) within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support,  $n = 69$ .

**Valentine (E6)**

*Message description:* The Valentine phish appeared to come from “eCard Delivery” with the email address, “do\_not\_reply@ecardalert.com.” The subject line said, “Happy Valentine’s Day! See who sent you an e-card. . .” There were three large red heart images at the top of the email. There was no personalization. The body of the email said, “A secret admirer wished you a Happy Valentine’s Day! Some of you may have heard about our employee greeting cards that can be used to acknowledge fellow employees. Click on the link below to view yours.” This was followed by a large link that said, “Your Card is Waiting,” and additional text that said, “If you are having trouble viewing the e-card please click [here](#).” “Would you like to send an e-card? Visit our [site](#). *Making someone’s day, one e-card at a time.* . . .” The email closed with, “This email may contain confidential and privileged information for the sole use of the intended recipient. If you are not the intended recipient, please contact the sender and delete all copies. Any review or distribution by others is strictly prohibited. Thank you.” The Valentine phish was sent January 22, 2018, prior to Valentine’s Day.

*Premise alignment (Method 1):* The alignment is categorized as Low—the primary premise related to Valentine’s Day does not align with a business process but is more personal in nature. However, the message mentions acknowledging a fellow employee, a sentiment that is recommended in the workplace. So, while the main premise

**Table 8:** valentine exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment Rating
Mimics a workplace process or practice	0 (Not applicable)
Has workplace relevance	2 (Low)
Aligns with other situations or events, including external to the workplace	6 (Significant)
Engenders concern over consequences for NOT clicking	2 (Low)
Has been the subject of targeted training, specific warnings, or other exposure	0 (Not applicable)
Total	10

**Table 9:** order confirmation exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	4 (Moderate)
Has workplace relevance	2 (Low)
Aligns with other situations or events, including external to the workplace	6 (Significant)
Engenders concern over consequences for NOT clicking	4 (Moderate)
Has been the subject of targeted training, specific warnings, or other exposure	0 (Not applicable)
Total	16

does not align with a business practice, it does play on the reader’s curiosity, aligning with the upcoming occasion of Valentine’s Day, of which most people are aware.

*Premise alignment (Method 2):* The overall alignment is 10 out of a possible 32. Table 8 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* All staff at NIST with an email address were targeted, from the human resources department, to finance, to bench scientists, to administrative support and all levels of management,  $n = 4\ 977$ .

**Order confirmation (E7)**

*Message description:* The order confirmation phish appeared to be from, “Order Confirmation” (auto-confirm@discontcomputers.com). Note the misspelling of “discount” in the email address. The subject line was personalized and said, “[Firstname Lastname]Your order has been processed,” with a space missing between the user’s last name and the word “Your.” At the top of the email was an image of several holiday packages, with the words, “Order Confirmation” in bold immediately below the holiday package image. There was no personalization in the body of the email, nor was there a greeting of any type. The email body text said, “Thank you for ordering with us. Your order has been processed. We’ll send a confirmation e-mail when your item ships.” This was followed by the words, “Order Details” in orange with, “Order: #SGH-2548883-2619437” (the order number was in blue text). The next section of the email said, “Estimated Delivery Date: 12/02/2016” (the date was in green text), “Subtotal: \$59.97,” “Estimated Tax: \$4.05,” and “Order Total: \$64.02” in bold. There was a large yellow button labeled with the text, “Manage order.” The button was followed by the text, “Thank you for your order. We hope you return soon for more amazing deals.” Near the bottom of the email was an image of a holiday snow globe and the text, “Need it in time for the holidays? Order before December 23 for free over-night shipping.” (“December 23” was in blue). Much smaller gray text below that said, “Unless otherwise stated, items sold are subject to sales tax in accordance with local laws. For more information, please view tax information” (“tax information” was in blue). Note the repeated word “in in,” a subtle mistake that is very difficult for users to notice, especially given the small gray font. Finally, at the very bottom of the email appeared three additional links, all in blue on a single line: “Return Policy | Privacy | Account.”

*Premise alignment (Method 1):* The alignment is categorized as Medium—the premise aligned for those who had purchasing authority in the OU and for those who had recently placed an order, a small subset of the whole OU. However, the training exercise took

**Table 10:** security token exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	2 (Low)
Has workplace relevance	6 (Significant)
Aligns with other situations or events, including external to the workplace	2 (Low)
Engenders concern over consequences for NOT clicking	2 (Low)
Has been the subject of targeted training, specific warnings, or other exposure	0 (Not applicable)
Total	12

place in December, when many people make on-line purchases for the holidays.

*Premise alignment (Method 2):* The overall alignment is 16 out of a possible 32. Table 9 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One OU within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support,  $n = 66$ .

#### Security token (E8)

*Message description:* The security token phish appeared to be from “Alerts” with the email address, “alerts@verifytoken.com.” The subject line was, “Verify Your Security Token Was Not Compromised.” The email was personalized using the format, “Lastname, Firstname, Middle Initial (Fed).” The body said, “Recently we have been made aware of a security breach in our security token product. Some of the tokens have been compromised and may need to be replaced. In order to find out if you [sic] token has been compromised, [Validate Your Security Token Here](#).” (Note the “you token” instead of “your token” here.) The email was signed “Rivest Shamir Adleman, Director of Identity and Access Management.” The email closed with smaller text that said, “This email may contain confidential and privileged information for the sole use of the intended recipient. Any review or distribution by others is strictly prohibited. If you are not the intended recipient, please contact the sender and delete all copies. Thank you.” If someone clicked the “Validate Your Security Token Here” link, they were taken to a data entry webpage, with the URL “secure.verifytoken.com.” The top of the webpage said, Token Security with a red background, followed by “Attention!!! Recently the safety of some security tokens has been compromised. Enter your username and six-digit number that is generated every 60 seconds by your security token and we will know if you will need a new token. Should you need a new token, you will be given contact information to request a new token, which will be shipped to you overnight.” This was followed by the text, “Account Login,” with fields labeled, “User ID” and “Password or Passcode.” There was a blue button labeled, “Login” and an “I’m not a robot” checkbox. At the bottom of the webpage was the text, “A passcode contains a PIN and a number from a security token.”

*Premise alignment (Method 1):* The alignment is categorized as Medium—the premise does not align at all for those personnel who do not have a security token (roughly 43% of all staff at the organization). Further, the premise does not align for those staff who expect any token checking and replacement would be conducted via the organization rather than a third party—likely a significant

**Table 11:** gift certificate exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	0 (Not applicable)
Has workplace relevance	0 (Not applicable)
Aligns with other situations or events, including external to the workplace	0 (Not applicable)
Engenders concern over consequences for NOT clicking	2 (Low)
Has been the subject of targeted training, specific warnings, or other exposure	0 (Not applicable)
Total	2

**Table 12:** adobe update exercise, premise alignment, method 2 ratings

Premise alignment element	Alignment rating
Mimics a workplace process or practice	0 (Not applicable)
Has workplace relevance	4 (Moderate)
Aligns with other situations or events, including external to the workplace	0 (Not applicable)
Engenders concern over consequences for NOT clicking	2 (Low)
Has been the subject of targeted training, specific warnings, or other exposure	-2 (Low)
Total	4

portion of the remaining 57% as the organization has a very strong posture regarding IT security.

*Premise alignment (Method 2):* The overall alignment is 12 out of a possible 32. Table 10 shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* All staff at NIST with an email address were targeted, from the human resources department, to finance, to bench scientists, to administrative support and all levels of management,  $n = 5024$ .

#### Gift certificate (E9)

*Message description:* The gift certificate phish appeared to be from “HR Rewards” with the email address, “certificates@great-restaurant-deals.com.” The subject line was “Your Restaurant Gift Certificate is here!” The body of the email contained a graphic in green, tan and gray showing people in a restaurant setting; the graphic was entitled, “Your Restaurant Gift Certificate is Attached!” The email was personalized using the format “Hi [Firstname Lastname]!” Following the greeting, the email read, “Your FREE complementary Restaurant Gift Certificate has arrived!” Note the misuse of complementary rather than complimentary. To the left of the message was the text, “Simply download and print the attached coupon and redeem it at any location or your choice! (Please be sure that the attachment’s barcode prints clearly.)” To the right of this text was a large tan rectangle containing the text, “25% savings!” Below this was the exclamation, “Happy dining!” Another rectangle in green contained the text, “Please see additional details and restrictions at the bottom of the official coupon, attached. Offer expires in 14 days from the date of

**Table 13:** phishing exercise data

Exercise number	Exercise name	Attack type	Number of cues	Premise alignment (Method 1)	Premise alignment (Method 2)	Difficulty rating	Actual phishing click rate
E1	Safety requirements ( $n = 73$ )	Link	7 (Few)	High	30 (High)	Very difficult	49.3% (36/73)
E2	Weblogs ( $n = 64$ )	Link	14 (Some)	High	24 (High)	Very difficult	43.8% (28/64)
E3	Unpaid invoice ( $n = 73$ )	Attachment	8 (Few)	High	24 (High)	Very difficult	20.5% (15/73)
E4	Scanned file ( $n = 62$ )	Attachment	6 (Few)	High	18 (High)	Very difficult	19.4% (12/62)
E5	New voicemail ( $n = 69$ )	Link	11 (Some)	Medium	14 (Medium)	Moderately difficult	11.6% (8/69)
E6	Valentine ( $n = 4\ 097$ )	Link	13 (Some)	Low	10 (Low)	Moderately/ Least difficult	11.0% (549/4 977)
E7	Order confirmation ( $n = 66$ )	Link	18 (Many)	Medium	16 (Medium)	Moderately difficult	9.1% (6/66)
E8	Security token ( $n = 5\ 024$ )	Credential harvesting via link	12 (Some)	Medium	12 (Medium)	Moderately difficult	8.7% (439/5 024)
E9	Gift certificate ( $n = 63$ )	Attachment	11 (Some)	Low	2 (Low)	Moderately/ Least difficult	4.8% (3/63)
E10	Adobe update ( $n = 63$ )	Link	12 (Some)	Low	4 (Low)	Moderately/ Least difficult	3.2% (2/63)

this email.” The message concluded with the text, “© 2014, All Rights Reserved.”

*Premise alignment (Method 1):* The alignment is categorized as Low—NIST’s human resources department does not award gift certificates. The premise capitalizes on the reader’s potential desire to for a 25% gift certificate to a restaurant of their choice.

*Premise alignment (Method 2):* The overall alignment is 2 out of a possible 32. [Table 11](#) shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One OU within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support,  $n = 63$ .

#### Adobe update (E10)

*Message description:* The Adobe update phish appeared to be from “Adobe Alerts” with the email address “adobeupdates@applt.net.” The subject line was, “Latest Adobe Security Update.” The email had a “High” importance level in the header information. The email was not personalized and no salutation was present. The first two lines of the email body contained “Security bulletin” and “Security updates available for Adobe Reader and Acrobat.” A “Release Date” of “February 12, 2015” was given, five days prior to the email sent date. The next line of text contained the following: “Vulnerability identifiers: APSB11-24.” This was followed by “CVE numbers: CVE-2011-1353, CVE-2011-2431, CVE-2011-2432, CVE-2011-2433, CVE-2011-2434, CVE-2011-2435, CVE-2011-2436, CVE-2011-2437, CVE-2011-2438, CVE-2011-2439, CVE-2011-2440, CVE-2011-2441, CVE-2011-2442.” The CVE numbers were followed by, “Platform: All.” Next was the single word, “SUMMARY” and the following text: “Critical vulnerabilities have been identified in Adobe Reader XI (11.0.04) and earlier versions for Windows and Macintosh, Adobe Reader 9.5.5 and earlier versions for UNIX, and Adobe Acrobat XI and earlier versions for Windows and Macintosh. These vulnerabilities could cause the application to crash and potentially allow an attacker to take control of the affected system.” The ‘summary’ was followed by the text, “Click here to download and install the update.” The last line of the email contained the text, “More details and solution information can be found [here](#).” There was no closing or contact information.

*Premise alignment (Method 1):* The alignment is categorized as Low—the premise aligns with the fact that nearly all staff had an Adobe Reader or Acrobat on their computers; however, the email-initiated update process did not align at all with the NIST-wide process for updating software.

*Premise alignment (Method 2):* The overall alignment is 4 out of a possible 32. [Table 12](#) shows the assigned rating for each element in the formulaic method and the sum.

*Target audience:* One OU within NIST, handling financial matters (ordering and invoice reconciliation), administrative program support, and technical program support.  $n = 63$

## Determining difficulty ratings

### Applying the Phish Scale

As described previously, the difficulty rating for an individual phishing message is determined first by categorizing the number of objectively observed cues and the premise alignment. This pair of categorizations is used to select the detection difficulty rating on the Phish Scale, shown in the conceptual framework in [Table 1](#), for the phishing message for an identified target audience. In [Table 13](#) the Phish Scale ratings are shown for each of the ten phishing exercises described in the last section, including the number of cues for each email (detail provided in [Supplementary Appendix A](#)), the premise alignment—Methods 1 and 2—from the exercise description and discussion with the training implementer), the difficulty rating (from the conceptual framework in [Table 1](#)), and the actual click rates for each exercise.

The table in [Supplementary Appendix A](#) contains the counts for each cue and a total count for each exercise. When counting cues in a given email message during analysis, it is important to note that these cue counts are based on our extremely careful scrutiny of the email messages; most email users are not going to notice or attend to all the available cues.

Note that in order to calculate the difficulty rating, the number of cues must be further categorized into *Few*, *Some*, or *Many*, in order of decreasing difficulty. Although some cues are more salient than others, we anticipate this is a reasonable first approximation. In this initial version of the Phish Scale, we propose the associated ranges as follows: the category labeled *Few* is represented by 1–8 cues, the category labeled *Some* by 9–14 cues, and the category labeled *Many* by 15 or more cues.

Further, we propose the premise alignment ratings are associated the following scores using Method 2: *High* alignment is obtained with an overall score of 18 and above, *Medium* alignment is obtained with an overall score between 10 and 18, and *Low* alignment is obtained with an overall score of 10 and below.

These ranges are based on our existing dataset; at this stage of scale development, the click rates inform the categorization of the cue counts and premise alignment scores. We fully expect these

**Table 14:** the Phish Scale—category coverage

Number of cues	Premise alignment	Detection difficulty	Exercise with click rate
Few (more difficult)	High	Very difficult	Safety requirements (49.3%), Unpaid invoice (20.5%), Scanned file (19.4%)
	Medium	Very difficult	–
	Low	Moderately difficult	–
Some	High	Very difficult	Weblogs (43.8%)
	Medium	Moderately difficult	New voicemail (11.6%) Security token (8.7%)
	Low	Moderately to least difficult	Valentine (11.0 %) Gift certificate (4.8%) Adobe update (3.2%)
Many (less difficult)	High	Moderately difficult	–
	Medium	Moderately difficult	Order confirmation (9.1%)
	Low	Least difficult	–

ranges may change with broader application of the Phish Scale to a larger variety of phishing emails. A larger corpus of phishing emails is needed to first inform the ranges and then an additional corpus of phishing emails is needed to validate the ranges.

It should be emphasized that the number of cues or the premise alignment alone does not determine the detection difficulty for a target audience; it is only when these elements are considered together for a target audience that a detection difficulty rating can be computed.

The Security token phish was the only exercise with a data entry component—after clicking the link users were taken to a webpage requesting their credentials. We report the data entry rates here rather than in Table 13. For the Security token phish, 24.4% (107/439) of clickers entered data on the credential-harvesting webpage. However, this is only 2.1% (107/5024) of the total number of employees who received the phishing email. Given that roughly 75% of clickers did not enter data on the webpage, it seems that additional suspicion was triggered on this page, likely due to being asked for credentials. This is certainly in line with the fact that mandatory yearly security awareness training at NIST in the past has focused heavily on not sharing credentials.

### Observations

Through these ten phishing exercises, we applied our Phish Scale to a variety of phishing attack types. This includes link-based attacks (Safety requirements, Weblogs, New voicemail, Valentine, Order confirmation, and Adobe update), attachment attacks (Unpaid invoice, which mimicked the real-world Locky ransomware attack, Scanned File, and Gift certificate), and a data entry or credential-harvesting attack (Security token). Now that we have used the Phish Scale to determine the detection difficulty rating for ten phishing exercises, there are a few observations we can make.

All of the exercises having a detection difficulty rating of *Very difficult* also have relatively high click rates (Safety requirements: 49.3%, Weblogs: 43.8%, Unpaid invoice: 20.5%, and Scanned file: 19.4%). However, the Weblogs exercise has many more cues than the other three exercises, and at 14 cues, was at the extreme end of the *Some* cues range (9 to 14).

All of the exercises having a detection difficulty rating of *Moderately difficult* have relatively lower click rates (New voicemail: 11.6%, Order confirmation: 9.1%, and Security token: 8.7%). Likewise, the exercises having the detection difficulty range *Moderately difficult* to *Least difficult* (Valentine: 11.0%, Gift certificate: 4.7%, and Adobe update 3.2%) have mid to lower click

rates. The Valentine and Security token exercises have a relatively larger and more varied sample than the other exercises which likely makes it more difficult to categorize the premise alignment. And finally, we do not have an exercise with a detection difficulty rating of *Least difficult*, calling attention to the need to apply the Phish Scale to additional exercises. In Table 14, we return to the conceptual framework given in Table 1 and note which categorizations are represented with operational data presented here.

### Limitations

This work is an early effort to characterize phishing message detection difficulty for email users situated in their normal email processing environments. As such, we acknowledge there are certainly limitations with this work at this time.

Current notable limitations in this work include: (i) the list of cues is long but not exhaustive; (ii) the uneven saliency of cues is not addressed in cue counts, but may be reflected indirectly in premise alignment; (iii) more experience is needed with categorizing premise alignment to improve assessment guidance; (iv) cue count ranges and premise alignment ranges (Method 2) need to be informed by additional data; and, (v) additional data are needed for scale validation. Further, the additional data used to inform the range components and ultimately the scale validation should be from diverse populations and sectors.

We anticipate that each of these limitations will be addressed as the Phish Scale is developed further.

## Discussion and future directions

### Click rates alone are insufficient: Why phishing detection difficulty matters

CISOs responsible for overseeing embedded phishing awareness training are often concerned when they observe click rates that are higher than expected. They are left wondering why click rates continue to be variable—possibly including large spikes—despite spending a significant amount of money and time training staff. CISOs must justify their cyber awareness training budgets and show a good ROI, lest their funding for such training be reduced. Unfortunately, if click rates continue to be high or variable, it is often—and we posit, incorrectly—perceived as due to ineffective training. We argue that this perception is fundamentally incorrect and hope to begin dispelling this perception through our development of a Phish Scale. Furthermore, we argue against focusing solely on phishing exercise

click rates, and instead strongly encourage the inclusion of reporting rates and reporting times as well; these metrics must be considered in conjunction, not in isolation, as early reporting can greatly improve mitigation efforts. Are reporting rates higher than click rates? Is time to first report sooner than time to first click?

We hope to frame the discussion around high click rates in a way that makes sense to CISOs and argue that high click rates can indicate that users are being exposed to new, difficult, and contextually relevant phishing campaigns. As part of a comprehensive phishing awareness cybersecurity program, we firmly believe difficult exercises actually improve user training effectiveness and awareness for real-world threats more than solely repeating the same or very similar, easier-to-detect phish. Click rates must be considered in conjunction with a deeper understanding of the phishing emails themselves and in light of reporting behavior as well. To this end, we have developed a Phish Scale to aid CISOs in better understanding and characterizing the detection difficulty of a given phishing exercise. Using operational data, the scale provides an indication of the difficulty email users in a target population will have detecting a particular phishing message. The Phish Scale addresses multiple components of phishing detection difficulty: cues, such as refs [15] and [16], and user context alignment [3].

Although our Phish Scale cue list is quite extensive, it is not exhaustive; however, phishing message developers are continually refining their methods and the cue list is easily extended. In moving towards a more formulaic approach to premise alignment categorization (Method 2) than presented in ref. [18], we now incorporate measures of perceived consequence severity and training effects. Greene *et al.* [3] found that clickers were concerned over consequences arising from *not* clicking, such as failing to be responsive to their job duties. In contrast, non-clickers were more concerned over consequences due to clicking, such as accidentally downloading malware. Additionally, concern over consequences varied depending on the premise of the phish. For example, it is likely that concern over consequences was much higher for the Weblogs exercise, with its implied consequence of disciplinary action, to include dismissal. We believe the addition of the more formulaic approach to premise alignment categorization that incorporates the elements: (i) *mimics workplace practice or practice*, (ii) *has workplace relevance*, (iii) *aligns with other situations or events*, (iv) *engenders concern over not clicking*, and (v) *targeted training, specific warnings, and other exposure*, make it easier for CISOs and training implementors to categorize premise alignment.

We expect that the three detection difficulty ratings we identified, *Very difficult*, *Moderately difficult*, and *Least difficult*, will eventually equate to validated click rate ranges. In speaking with CISOs, we anticipate ranges roughly along these lines: the *Very difficult* category having click rates near or above 20%, the *Moderately difficult* category having click rates in the approximately 10 to 20% range, and the *Least difficult* category having click rates below 10%. We plan to inform the actual ranges with additional empirical data; the ten exercises presented here are a start.

It is early days for the Phish Scale, however, we believe the conceptual framework has promise when we consider the projected detection difficulty rating and the actual click rates for the ten exercises we examined. Additionally, we stress the Phish Scale components are still in development. We know all cues do not have equal salience. Finding an abbreviated method for CISOs to characterize premise alignment has proven difficult. That said, we believe there is value in bucketing phishing message detection difficulty to tailor training and contextualize click rates for both training and actual threats.

## Differential cue salience: Not all cues are created equally

Capturing the effect of phishing message cues is difficult, as not all cues are created equally. The saliency and effect of any particular phishing cue varies, determining whether it is perceived as a suspicion indicator versus a compelling hook. This aspect of phishing message characteristics is important to note. Whether a cue is perceived as a phish indicator versus a hook depends on the user and the user's context when processing the email. Well-known phishing indicators such as misspellings and grammar errors are often regarded by email users as suspicion-generating, and when noticed can lead to additional user scrutiny of the message for more phishing indicators. Another undisputed phishing characteristic is urgency. Its use is so common that it should be a red flag; however, urgency is legitimately common-place in today's world, diluting its suspicion-generating signal strength. Additionally, urgency inhibits slower, more deliberate (System 2 [12]) processing, making it more a hook enhancer than a red flag.

We suspect that saliency for some cues may actually be reflected in the premise alignment categorization. For example, in the Unpaid Invoice exercise, the sender's designation as a (Fed) in the message's from field served as a compelling cue for those who clicked as it mimicked a workplace convention and practice. The presence of this message element strengthened the overall premise alignment and in doing so nudged the premise alignment categorization higher.

## Categorizing premise alignment with user context

In this early version of the Phish Scale, we use the terms *High*, *Medium*, and *Low* to bucket premise alignment for a target audience into intuitive high-level categories with associated definitions (Method 1). A more formulaic approach is documented in this work: premise alignment categorization (Method 2). While these methods are sufficient for the beginning phase of scale development, we may seek to refine the characterization methods for the categorical variables in future work, by investigating contextual relevance measures and scales. "Contextual" is a part of existing scales in other domains such as, "A Contextual Measure of Achievement Motivation" [40] and "contextual performance" as a dimension of individual work performance [41]. How might such existing scales and measures be leveraged for use in the phishing domain? Additionally, how do we account for changes in context over time?

Changes in contextual relevance may occur over quite long timescales, as someone slowly adds or changes job responsibilities over the years of their career, or very short timescales, as some event that day/week/month may trigger heightened contextual relevance. For example, Greene *et al.* [3] explained that users were concerned over a real-world vendor invoice that was unpaid, leading to temporarily heightened contextual relevance for the unpaid invoice phishing email. Daily events, such as expecting or missing a phone call, can temporarily heighten the contextual relevance of a "new voicemail" phishing email. Factors such as being busy, stressed, or rushed can also fluctuate widely during a work day. It is likely the case that there is a relatively fixed component of user context, in addition to a more time-sensitive, variable component. The current Phish Scale does not break down context and associated premise alignment into these subcomponents. It is unclear whether such a fine-grained distinction is indeed necessary at this point.

Although it may be quite feasible to discern premise alignment with finer granularity than our existing categories, this may actually be superfluous for the intended audience of the Phish Scale. With our goal of developing a simple, easy to use Phish Scale for CISOs

and those responsible for implementing and overseeing phishing awareness training programs, it is likely the case that *High*, *Medium*, and *Low* categories for premise alignment are sufficient. The important point we seek to emphasize with our Phish Scale is that a highly relevant context makes it extremely difficult for users to detect phishing emails. The greater the contextual relevance, the less likely a user is to notice, attend to, and think deeply about suspicious email cues. Daily stressors such as time pressure in general reduce the cognitive resources that users have available to dedicate to email processing. When cognitive resources are reduced, it makes it more likely that users will engage in faster, heuristic, System 1 processing rather than thoughtful, slower, deeper System 2 processing [12].

A final point with respect to user context and premise alignment has to do with the target audience size: categorizing premise alignment becomes more difficult as the size of the target audience increases. With a larger target audience, there is typically a much greater variety of work responsibilities present and a wider variety of user contexts, which may or may not align with a phishing email premise. Depending on the phishing message premise, there may be benefit in additional review of click rates for groups based on work roles rather than for a large target audience with differing roles. In any case, premise alignment must be determined by someone who has a good grasp of the collective context of work for the specific target audience of a phishing message. To have meaning, it cannot be otherwise.

### Comparing phishing data across sectors

Although cross-exercise and cross-sector phishing comparisons are frequently made, and are indeed quite valuable, interpretation of such comparisons still pose significant challenges. In particular, when the level of phishing detection difficulty can vary so dramatically based on user context and premise alignment, it is in some sense a meaningless comparison without a basic understanding and assessment of: (i) characteristics of the phishing email itself and (ii) characteristics of the target user population. More specifically, one must understand the premise and cues contained within a given phish in conjunction with the work context of the target user population. Toward this end, we believe our Phish Scale shows great promise as a tool to help frame data sharing on click rates and reporting rates across exercises, organizations, and sectors.

As we refine and mature this tool with input from the larger usable security community, we hope to move the Phish Scale out of the research community and into operational use. For instance, we believe that beyond providing benefits to CISOs and phishing training implementers, our Phish Scale could also provide significant value to joint organizations responsible for sharing cyber threat intelligence data. For example, the Federal Bureau of Investigation (FBI), has an InfraGard program, a partnership between the FBI and the private sector dedicated to sharing information and intelligence [42]. There are other such collaborative programs as well, for example, the National Cyber-Forensics and Training Alliance (NCFTA) which is a nonprofit partnership between private industry, government, and academia working together to disrupt cybercrime [43]. Phishing in particular, and social engineering in general, are active threats across all industry verticals. By providing a phishing difficulty rating framework, our Phish Scale can help facilitate collaboration using a common language surrounding human phishing threat detection.

### Future work

We encourage other usable security researchers and practitioners to use our Phish Scale, apply it to a much wider variety of phishing emails, and test its predictions against both existing phishing training exercise data, and ultimately against real-world phishing emails as well. We plan to continue applying our Phish Scale to a larger corpus of additional emails for which we have click rate and premise alignment data, and plan to partner with external entities to do the same. Unfortunately, our access to concurrent reporting data is more limited. A notable challenge of conducting research with operational workplace data is that there is often a tradeoff between experimental control and ecological validity. In this case, we had the benefit of extremely high ecological validity, as users were in their normal workplace settings with their normal tasks and email loads, but without the control necessary to capture reporting rates at the time of the phishing exercises. Nonetheless, the benefit of having new, *in situ* workplace data for approximately 5000 employees offers an important contribution to the phishing literature and to the larger usable security community.

Beyond applying and testing the current Phish Scale with additional data, we intend to explore new scale components as well. We would like to investigate incorporating work on personality factors, curiosity, distractedness, concern for security, and ultimately folding the various components of our Phish Scale into a lens model, an application of multiple regression often used in judgement and decision-making research. This would build upon prior lens modeling work by Tamborello and Greene [44] and Molinaro and Bolton [15]. Additional modeling and simulation research could explore the predicted click rates and reporting rates for different combinations of cues, context alignment, personality types, and phishing premises. How do different combinations affect phishing susceptibility? For example, consider this combination: users scoring high on conscientiousness, with a financial work context, who receive a phishing email with an authoritarian/time-sensitive transfer of funds premise, and very few suspicious cues. What if everything were the same but the work context, is that difference alone sufficient for someone to catch this phish? While we believe that context may trump all, additional research is necessary to see in which scenarios this holds, as well as, how and when it may change. One could simulate—with a well-validated model—the large number of possible combinations, to determine where to focus research and training intervention efforts based on quantified predicted risk metrics, such as the likelihood of clicking versus reporting.

### Broader implications

The Phish Scale—and indeed phishing in general—is part of a much larger research agenda that addresses a spectrum of usable security issues. For instance, understanding risk, including human risk, is a key component of any organization's cybersecurity strategy, and risk management frameworks play an important role in helping maintain security and privacy [45]. Ultimately, we hope our Phish Scale can be used to help CISOs better understand and characterize their organization's phishing risk, by essentially profiling the types of phishing premises their users are more or less susceptible to as well as the organization's actual threats. Such data can be used to prioritize training efforts on more targeted interventions, and to prioritize investigative efforts for real-world suspected phishes. Targeted training interventions will likely need to move beyond embedded phishing exercises, especially for repeat clickers. In-person seminars, posters, informal lunch and learn sessions, and so on, are all part of a larger security awareness program. Additional

interventions may include special email graphical user interface (GUI) elements or flagging, or perhaps more aggressive email filtering for certain users or groups based on their risks and job responsibilities.

In addition to risk profiling and targeted training, future work is also needed to understand how new technological email security measures will impact phishing. In particular, government agencies are quickly moving toward email authentication by implementing protocols such as Domain-based Message Authentication, Reporting, and Conformance (DMARC) and Domain Keys Identified Mail (DKIM) per the Department of Homeland Security (DHS) Binding Operational Directive 18-01 [46]. How will that affect the phishing space? On the other hand, pretexting is already gaining in popularity and will likely continue to do so, especially if new technological solutions prevent or threaten the success of certain more “traditional” phishing email scams. As advances in technological protections make some attacks less effective, or even one day obsolete, the attacks will not stop, but rather will transition and evolve in response. For instance, it seems likely that other out-of-band social engineering methods will continue to gain in popularity. Phishing is but one component of a much larger social engineering problem facing the cybersecurity field. Future work should examine how lessons learned in the phishing domain may inform other varieties of social engineering problems as well.

## Supplementary data

Supplementary data is available at *Journal of Cybersecurity* online.

## Acknowledgements

The authors gratefully acknowledge NIST’s Information Technology Office and Networking Division and our Information Technology Security Officer partner for their support throughout this project, as well as, Susanne Furman, who created the phishing message images contained in this document. We also thank the CISOs and training implementers who spoke with us regarding their phishing awareness programs, as well as anonymous reviewers for their comments.

*Conflict of interest statement.* This effort was funded by the National Institute of Standards and Technology. Any mention of commercial products or reference to commercial organizations is for information only; it does not imply recommendation or endorsement by the National Institute of Standards and Technology nor does it imply that the products mentioned are necessarily the best available for the purpose.

## References

1. Cybersecurity Ventures. 2019 Cybercrime Report, <https://www.herjavecgroup.com/the-2019-official-annual-cybercrime-report/> (15 May, 2020, date last accessed).
2. Greene KK, Steves M, Theofanos M. No phishing beyond this point. In *IEEE Comp Cybertrust Column* 2018;51:86–89.
3. Greene KK, Steves M, Theofanos M *et al.* User context: an explanatory variable in phishing susceptibility. USEC NDSS 2018. Usable Security Workshop at the Network and Distributed Systems Security Symposium. 18 February 2018. San Diego, CA. DOI: <https://dx.doi.org/10.14722/usec.2018.23016>
4. Newman LH. What spammers could do with your hacked Facebook data. *Wired*, 2018, <https://www.wired.com/story/facebook-hack-data-spammers/> (May 2020, date last accessed).
5. Healthcare CyberGard Annual Conference, Charlotte, October 2018, <https://www.ncfta.net/healthcare-cybergard-annual-conference-charlotte-october-2018/> (May 2020, date last accessed).

6. Information Security and Privacy Advisory Board, <https://csrc.nist.gov/CSRC/media/Projects/ISPAB/documents/minutes/ispab-june-2018-meeting-minutes.pdf> (May 2020, date last accessed).
7. Sawyer BD, Hancock PA. Hacking the human: the prevalence paradox in cybersecurity. *Human Factors* 2018;60:597–609.
8. Levari DE, Gilbert DT, Wilson TD *et al.* Prevalence-induced concept change in human judgement. *Science* 2018;360:1465–7.
9. Rogers RW. A protection motivation theory of fear appeals and attitude change. *J Psychol* 1975;91:93–114.
10. Wang J, Li Y, Rao R. Coping responses in phishing detection: an investigation of antecedents and consequences. *Inf Syst Res* 2017;28:378–96.
11. Vishwanath A, Herath T, Chen R *et al.* Why do people get phished? Testing individual differences in phishing vulnerability within an integrated, information processing model. *Decis Supp Syst* 2011;51:576–86.
12. Kahneman D. *Thinking, Fast and Slow*. New York: Farrar, Straus and Giroux, 2011.
13. Vishwanath A, Harrison B, Ng YJ. Suspicion, cognition, and automaticity model of phishing susceptibility. *Comm Res* 2018;45:1146–66.
14. Williams EJ, Hinds J, Joinson AN. Exploring susceptibility to phishing in the workplace. *Int J Human-Comp Stud* 2018;120:1–13.
15. Molinaro KA, Bolton ML. Evaluating the applicability of the double system lens model to the analysis of phishing email judgments. *Comp Sec* 2018;77:128–37.
16. Parsons K, Butavicius M, Pattinson M *et al.* Do users focus on the correct cues to differentiate between phishing and genuine emails?. *Australasian Conference on Information Systems* 2015; <https://arxiv.org/abs/1605.04717> (May 2020, date last accessed).
17. Caputo D, Pfleeger SL, Freeman J *et al.* Going spear phishing: exploring embedded training and awareness. *IEEE Sec Priv* 2014;12:28–38.
18. Steves M, Greene K, Theofanos M. A phish scale: rating human phishing message detection difficulty. USEC NDSS 2019. *Usable Security Workshop at the Network and Distributed Systems Security Symposium*. February, 2019. San Diego, CA. DOI: <https://dx.doi.org/10.14722/usec.2019.23028>.
19. Blythe, M Petrie, H Clark JA. F for fake: Four studies on how we fall for phish. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, May 2011, Vancouver, Canada, ACM. 2011, pp. 3469–78.
20. Canfield CI, Fischhoff B, Davis A. Quantifying phishing susceptibility for detection and behavior decisions. *Human Fact* 2016;58:1158–72.
21. Downs JS, Holbrook M, Cranor LF. *Decision strategies and susceptibility to phishing*. In: *Proceedings of the Second Symposium on Usable Privacy and Security (SOUPS '06)*, ACM, 2006, pp. 79–90.
22. Furnell S. Phishing: can we spot the signs? *Comp Fraud Sec* 2007;2007:10–15.
23. Grazioli S. Where did they go wrong? An analysis of the failure of knowledgeable internet consumers to detect deception over the internet. *Group Decis Negot* 2004;13:149–72.
24. Hadnagy C, Fincher M. *Phishing Dark Waters*. Hoboken, NJ: Wiley, 2015.
25. Jakobsson M. The human factor in phishing. *Priv Sec Consum Inform* 2007;7:1–19.
26. Jakobsson M, Finn P. Designing and conducting phishing experiments. In: *IEEE Technology and Society Magazine, Special Issue on Usability and Security*, IEEE, 2007.
27. Karakasioti A, Furnell S, Papadaki M. Assessing end-user awareness of social engineering and phishing. In: *Australian Information Warfare and Security Conference*, School of Computer and Information Science, 2006, Edith Cowan University, Perth, Western Australia.
28. Parsons K, McCormac A, Pattinson M *et al.* Phishing for the truth: a scenario-based experiment of users behavioural response to emails. In: *IFIP International Information Security Conference*, Springer, Berlin, Heidelberg, 2013, pp. 366–78.
29. Wang J, Herath T, Chen R *et al.* Phishing susceptibility: an investigation into the processing of a targeted spear phishing email. *IEEE Trans Prof Comm* 2012;55:345–62.
30. Wright R, Chakraborty S, Basoglu A *et al.* ‘Where did they go right?’ Understanding the deception in phishing communications. *Group Decis Negot* 2010; 19:391–416.

31. Han Y, Shen Y. Accurate spear phishing campaign attribution and early detection. In: *Proceedings of the 31st Annual ACM Symposium on Applied Computing*, ACM, 2016, pp. 2079–86.
32. Kim D, Kim JH. Understanding persuasive elements in phishing e-mails: a categorical content and semantic network analysis. *Online Inform Rev* 2013;37:835–50.
33. Egelman S, Cranor LF, Hong J. You've been warned: An empirical study of the effectiveness of web browser phishing warnings. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, ACM, 2008, pp. 1065–74.
34. Tsow A, Jakobsson M. Deceit and deception: a large user study of phishing. Indiana University, Technical report TR649, 2007.
35. Dhamija R, Tygar JD, Hearst M. Why phishing works. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, ACM, 2006, pp. 581–90.
36. Fogg BJ. Prominence-interpretation theory: explaining how people assess credibility online. In *CHI'03 Extended Abstracts on Human Factors in Computing Systems*, ACM, 2003, pp. 722–23.
37. Alsharnouby M, Alaca F, Chiasson S. Why phishing still works: user strategies for combating phishing attacks. *Int J Human-Comp Stud* 2015;82:69–82.
38. Perkins D, Salomon G. Transfer of learning. *Int Encycl Edu* 1992;2:6452–57.
39. Ransomware, <https://www.trendmicro.com/vinfo/us/security/definition/RANSOMWARE> (May 2020, date last accessed).
40. Smith RL. A contextual measure of achievement motivation: Significance for research in counseling. *VISITAS Online*, 2015.
41. Koopmans L, Bernaards CM, Hildebrandt VH *et al.* Measuring individual work performance: identifying and selecting indicators. *Work* 2014;48:229–38.
42. FBI, Federal Bureau of Investigations. <https://www.fbi.gov/about/partnerships/infragard> (May 2020, date last accessed).
43. NCFETA, National Cyber-Forensics and Training Alliance, <https://www.ncfta.net>. (May 2020, date last accessed).
44. Tamborello FP, Greene KK. Exploratory lens model of decision-making in a potential phishing attack scenario. National Institute of Standards and Technology Interagency Report, NISTIR 8194, October 2017. <https://doi.org/10.6028/NIST.IR.8194>.
45. NIST, National Institute of Standards and Technology, 'Special Publication 800-37, Revision 2, Risk Management Framework for Information Systems and Organizations—A System Life Cycle Approach for Security and Privacy', 2018.
46. DHS, Department of Homeland Security. Binding Operational Directive BOD-18-01, 2017, <https://cyber.dhs.gov/assets/report/bod-18-01.pdf> (May 2020, date last accessed).