

# Anti-Phishing Training Does Not Work: A Large-Scale Empirical Assessment of Multi-Modal Training Grounded in the NIST Phish Scale

Andrew T. Rozema  
Purdue University  
arozema@purdue.edu

James C. Davis  
Purdue University  
davisjam@purdue.edu

**Abstract**—Social engineering attacks using email, commonly known as phishing, are a critical cybersecurity threat. Phishing attacks often lead to operational incidents and data breaches. As a result, many organizations allocate a substantial portion of their cybersecurity budgets to phishing awareness training, driven in part by compliance requirements. However, the effectiveness of this training remains in dispute. Empirical evidence of training (in)effectiveness is essential for evidence-based cybersecurity investment and policy development. Despite recent measurement studies, two critical gaps remain in the literature: (1) we lack a validated measure of phishing lure difficulty, and (2) there are few comparisons of different types of training in real-world business settings.

To fill these gaps, we conducted a large-scale study ( $N = 12,511$ ) of phishing effectiveness at a US-based financial technology (“fintech”) firm. Our two-factor design compared the effect of treatments (lecture-based, interactive, and control groups) on subjects’ susceptibility to phishing lures of varying complexity (using the NIST Phish Scale). The NIST Phish Scale successfully predicted behavior (click rates: 7.0% easy to 15.0% hard emails,  $p < 0.001$ ), but training showed no significant main effects on clicks ( $p = 0.450$ ) or reporting ( $p = 0.417$ ). Effect sizes remained below 0.01, indicating little practical value in any of the phishing trainings we deployed. Our results add to the growing evidence that phishing training is ineffective, reinforcing the importance of phishing defense-in-depth and the merit of changes to processes and technology to reduce reliance on humans, as well as rebuking the training costs necessitated by regulatory requirements.

## 1. Introduction

Phishing attacks critically impact organizations through operational incidents, data breaches, and financial losses. They constitute the primary attack vector to breach organizations as highlighted by industry studies [1] and the US Federal Bureau of Investigation (FBI) [1], [2]. To prevent this class of attacks, organizations employ technical controls (e.g., spam filters, attachment analyzers) complemented by cybersecurity awareness training [3], [4].

Although these anti-phishing trainings are costly, their effectiveness is in question. Organizational security awareness programs often focus on compliance metrics rather than

meaningful behavioral change [5]. Meanwhile, recent large-scale evidence suggests that current approaches have limited impact [4]. To better evaluate the effectiveness of anti-phishing training, two questions arise. First, are some kinds of training more effective than others? Second, how does their effectiveness vary with the difficulty of the phishing lure?

In this paper, we report the results of a measurement study to address these questions. We situate the study in an operational environment [6], set at a US-based international fintech firm with **12,511 employees** across multiple business units. We compared three experimental conditions: a control group (no training), traditional lecture-based training, and lectures augmented with interactive phishing exercises. Subjects in each condition received lures of varying complexity, which we assigned difficulty levels using the NIST Phish Scale. This study design enables rigorous evaluation of training effectiveness across standardized difficulty levels on a large population, improving on previous works that had either been small-scale or lacking in systematic difficulty assessment [7], [8].

Our findings were unfavorable for the conventional anti-phishing trainings available from cybersecurity training vendors. Neither mode of training produced statistically significant improvements in phish click reduction (average: 10.4%) or phish reporting behavior (average: 9.6%). However, the NIST Phish Scale did predict user behavior, with phishing lure difficulty having a highly significant effect on click-through rates ( $F(2, 12086) = 41.415$ ,  $p < 0.001$ ,  $\eta^2 = 0.007$ ). Thus our study shows that, while the NIST Phish Scale provides a means of measuring the difficulty of different lures, the trainings we applied did not help subjects perform better on those lures.

We summarize our main contributions below:

- 1) **Large-Scale Validation of the NIST Phish Scale.** We provide the most comprehensive empirical validation of the NIST Phish Scale to date, demonstrating its predictive utility across 12,511 participants in an operational environment, showing strong correlation between difficulty ratings and user susceptibility.
- 2) **Realistic Assessment of Training Effectiveness.** Our results challenge optimistic assumptions about the effectiveness of phishing training. We found no statistically significant impact of training on either click rates or

reporting behavior. Although these findings align with recent large-scale studies from the healthcare [4] and banking sectors [9], our study addresses explicit calls from the research community for large-scale ecologically valid studies combining standardized assessments of difficulty inside of an operational context [6], [10].

**Significance:** Our findings have important implications for researchers, organizations, and regulators. For the research community, we provide empirical assessment of training effectiveness using rigorous methodology and standardized difficulty measures. For organizations, our results suggest the need for realistic expectations about training outcomes and highlight the importance of complementing awareness programs with robust technical controls rather than relying primarily on human-centered defenses. Our validation of the NIST Phish Scale’s predictive utility enables organizations to use this tool for benchmarking the difficulty of their phishing simulations. For regulators, our work contributes to the growing body of evidence that organizational cybersecurity strategy should emphasize technical defenses [1] rather than placing the burden on humans — requiring anti-phishing training is costly and has minimal effect.

## 2. Background and Related Work

In a phishing attack, attackers send messages (*phishing lures*) to targets whom they wish to persuade to act against their own interests. Attackers seek results such as sensitive information, account credentials, or the ability to control the victim’s machine [1], [11]. Modern phishing campaigns target individuals and organizations across multiple attack vectors, applying conventional email phishing as well as SMS-based attacks (“smishing”) and voice-based deception (“vishing”) [12], [13]. Recent attacks have begun to incorporate Generative AI tools to produce tailored and realistic messages. The severity, persistence, and increasing sophistication of these attacks has established phishing detection and prevention as a critical focus within both academic cybersecurity research and industry practice [6], [14].

In this section, we summarize current defenses against phishing (§2.1), discuss metrics for evaluating phishing effectiveness (§2.2), and review prior empirical studies (§2.3) and knowledge gaps (§2.4).

### 2.1. Phishing Defenses

Phishing poses substantial risks to organizations, who therefore erect multiple layers of defenses against these attacks. We distinguish automated defenses from human-focused ones.

**2.1.1. Automated Defenses.** Organizations employ a range of technological defenses to mitigate phishing risks. These include spam filters, attachment analyzers, and URL checkers [1], [14], [15]. More recently, machine learning approaches have been applied to detect phishing content. Long Short-Term Memory (LSTM) models and optimized

TABLE 1: NIST Phish Scale for Detection Difficulty [32], [33].

Cues \ Alignment	Weak	Medium	Strong
Few	Medium	Hard	Hard
Some	Easy-to-Medium	Medium	Hard
Many	Easy	Medium	Medium

Random Forest algorithms [16] have demonstrated high accuracy in detecting phishing attempts. As with everything in cybersecurity, phishing is a cat-and-mouse game. Adversaries have begun to leverage artificial intelligence to create more convincing lures. Large language models (LLMs) can generate personalized and context-aware phishing emails that are difficult for traditional filters to detect [3]. Techniques like typosquatting, domain spoofing, and watering hole attacks exacerbate this issue [1], [14].

**2.1.2. Human-Focused Defenses: Training and Awareness.** Defending against phishing is understood not as solely a technical problem, but rather as a socio-technical one [6], [13]. Spear phishing and emerging phishing attacks present sophisticated threats that overcome automated defenses [17], [18], so human awareness of phishing remains necessary to successful defense [19], [20].

There are many approaches to training humans to be resilient to phishing attacks. Methods range from lecture-style training involving passive videos and readings [21] to simulation feedback providing immediate feedback after interaction with phishing simulations [22], [23]. Interactive training uses gamification or scenario-based learning to engage participants [24], [25], while just-in-time training alerts users in real time during risky behavior [26].

Anti-phishing trainings are frequently included within the cybersecurity awareness trainings ubiquitous in organizations. Legal mandates drive the widespread adoption of such training programs, with HIPAA requiring security awareness training for healthcare organizations [27], GDPR encouraging staff training on cybersecurity practices [28], and PCI DSS mandating training for those handling cardholder data [29]. Industry standards and frameworks further reinforce these requirements, with ISO 27001 calling for human resource security controls including awareness programs [30] and the NIST Cybersecurity Framework advising phishing simulations under the “Identify” function [31].

### 2.2. Measuring Effectiveness of Phishing Training

The effectiveness of phishing training can be assessed based on how resilient a target becomes to phishing lures. §2.2.1 describes how lures can be categorized by their complexity, and then §2.2.2 gives the standard metrics for target resilience.

**2.2.1. Measuring Training Lure Difficulty.** To assess phishing lure difficulty, the NIST Phish Scale categorizes them along two dimensions as depicted in Table 1. *Phishing Cues* represent observable errors or inconsistencies in

TABLE 2: Summary of Prior Experiments on Phishing and Training Interventions. Studies are organized by research context (Lab vs. Real-World). *Context*: “Lab” denotes research simulations, where subjects are not in their work roles. “Real-World” denotes studies where participants were employees acting in their job capacity. *Complexity measure*: Whether the study controlled for the difficulty of the phishing tasks. ✓ indicates the use of an open control such as NIST Phish Scale, and ✗ indicates no control. *Hypotheses*: How these works informed the hypotheses in the present study (cf. Table 3).

Work	Context	Population	Sample size	Complexity Measure	Compared Training Approaches?	Hypotheses
[8] (2024)	Lab	University students	96	✗	✗	H3
[34] (2024)	Lab	University students	117	✓ (NIST Phish Scale)	✗	H4
[7] (2022)	Real-world	Employees	119	✗	✗	H2
[22] (2017)	Real-world	University students	150	✗	✓	H2, H3
[23] (2014)	Real-world	Employees	1,359	✗	✓	H2, H4
[9] (2024)	Real-world	Banking employees	“Thousands”	✗	✗	H2
[35] (2023)	Real-world	Finance employees	5,000	✗	✗	H2
[4] (2025)	Real-world	Healthcare employees	19,500+	✗	✗	H2
<b>Ours</b>	Real-world	Fintech org.	12,511	✓ (NIST Phish Scale)	✓	—

the email that might alert users to its fraudulent nature, such as spelling mistakes, suspicious URLs, or formatting irregularities [34]. *Premise Alignment* measures the relevance and familiarity of the email’s subject matter to the recipient’s typical email environment and organizational context [36]. A high alignment, messages that closely resemble emails users receive daily, combined with few obvious cues yields the most deceptive and challenging-to-detect phishing emails [17]. The scale allows organizations to benchmark lure complexity and contextualize simulation outcomes, though it requires sufficiently large sample sizes for reliable use and consistent application [37].

**2.2.2. Training Outcome Metrics.** The effectiveness of anti-phishing training is commonly measured in three ways:

- *Open Rate* captures the percentage of emails opened, though this measurement can be deceptive since different email clients and security settings can affect whether an email is recorded as opened [36] we therefore largely ignore this measure.
- *Click-Through Rate (CTR)* measures the percentage of recipients who click on malicious links, serving as a primary indicator of phishing susceptibility [22], [23].
- *Reporting Rate* tracks phishing incidents reported to security teams, representing users’ proactive security behaviors and ability to recognize suspicious content [18], [38].

## 2.3. Empirical Studies of Anti-Phishing Training

For reviews of knowledge about susceptibility to phishing and about anti-phishing training, we refer the reader to recent works by Marshall *et al.* [21] and Franz *et al.* [10]. This section focuses on the aspect of the literature to which this study contributes. We discuss the settings of prior work

(§2.3.1), the assessments of different trainings (§2.3.2), and measurement and evaluation methods (§2.3.3).

**2.3.1. Scale and Context of Prior Work.** Many phishing training studies have been conducted in controlled laboratory settings, where researchers benefit from a high degree of experimental control. Studies of this sort have provided numerous insights, including that active learning and immediate feedback can improve short-term detection rates, and have explored how individual differences affect training outcomes [7], [21]. However, as noted by Franz *et al.* [10], these controlled settings may lack ecological validity. In laboratory environments, participant samples are typically small and consist of homogeneous groups such as academics and students. Due to this, and the artificial nature of lab tasks, these studies may not accurately reflect the complexities and distractions of real-world environments. This raises questions about the generalizability of laboratory findings to operational contexts, where additional factors can alter training effectiveness.

In response, a growing number of studies have begun examining phishing training in real-world organizational environments. Doing *et al.* [9] analyzed phishing reporting behaviors across thousands of bank employees, while Ho *et al.* [4] conducted a comprehensive evaluation with over 19,500 healthcare employees across 10 simulated phishing campaigns, reporting minimal effectiveness of both annual cybersecurity awareness training and embedded anti-phishing training. To date, however, these larger studies have lacked a standardized method for assessing the difficulty of the phishing lures in question [4], [9], or have focused primarily on observational analysis rather than controlled training interventions with standardized difficulty measures [9].

Our study contributes to the evidence on real-world human phishing resilience, with data from approximately

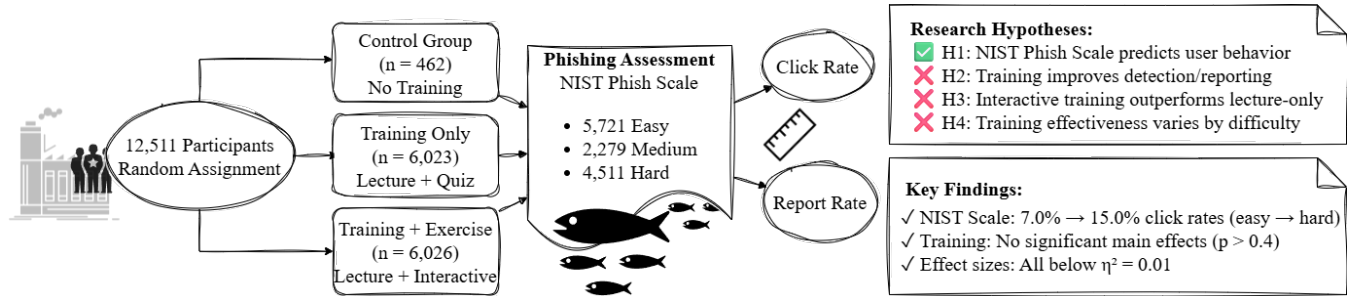


Figure 1: Illustration of our methodology. Our two-factor design randomly assigns different trainings to different subjects, and then randomly sends phishing lures of varying complexity to the subjects. We compared between-subject performance using lure click rate and lure report rate in order to test four hypotheses derived from prior work.

12,000 subjects in a fintech enterprise, while maintaining experimental control and measuring lure difficulty with the NIST Phish Scale framework. This approach allows us to examine whether the modest training effects we observe align with Ho *et al.*'s findings, while providing additional insights into how training effectiveness varies by phishing sophistication.

**2.3.2. Training Intervention Approaches.** Franz *et al.* [10] provide a comprehensive taxonomy identifying five categories of phishing interventions: training, awareness, education, notifications, and design interventions. Our study contributes empirical data within their training category, specifically comparing lecture-based delivery against interactive modalities. While Schöni *et al.* [8] demonstrated effectiveness of personalized training approaches, and our work examines scalable training modalities that can be deployed across demographically diverse cohorts without individualization, recent large-scale studies suggest that the practical benefits of such training interventions may be more limited than previously assumed [4]. Our study builds on this foundation by applying the NIST Phish Scale to measure lure difficulty in line with current guidance, operationalizing the scale at unprecedented scale, and measuring outcomes (clicking, reporting) post-intervention, allowing for practical comparison across conditions and lure complexity.

**2.3.3. Measurement and Evaluation Methods.** A critical gap in prior work has been the lack of standardized difficulty assessment for phishing simulations. Barrientos *et al.* [39] conducted early validation work on the NIST Phish Scale, while Dawkins and Jacobs [33] provided the official implementation guidelines. However, no studies have operationalized the scale at enterprise level nor examined how difficulty interacts with training effectiveness. We applied the scale across varied training modalities and measuring multiple outcome measures (clicking, reporting) post-intervention.

## 2.4. Research Gaps

Table 2 summarizes the preceding analysis as two gaps. First, prior work has primarily focused on individual-level factors influencing phishing susceptibility. The interaction

between training modality and phishing difficulty — a key question for optimizing organizational defenses — has received little attention. Second, despite widespread adoption of phishing training, few studies rigorously assess outcomes against standardized difficulty measures or operate at sufficient scale for robust inference.

This study addresses those gaps by applying the NIST Phish Scale [33] to evaluate training effectiveness across varied phishing difficulty levels in an enterprise deployment with over 12,000 participants. In this study, our objectives are to: (1) operationalize the NIST framework at scale, (2) examine how difficulty influences detection and reporting across training modalities, (3) provide comparative evidence on lecture-based versus interactive approaches, and (4) offer guidance for large-scale awareness efforts under conditions of modest effectiveness [4]. This enables evidence-based decision-making on phishing defense strategies that account for both training design and message complexity.

## 3. Methodology

In this section, we describe the methodology of our measurement study, including our hypotheses, experimental design, participant demographics, training modalities, and phishing simulations. An illustration is given in Figure 1. Our two-factor study design controls for the effect of training modality and task complexity. We make the first real-world application of the NIST Phish Scale framework [33], allowing us to evaluate the effectiveness of different training modalities across standardized difficulty levels in an operational environment.

### 3.1. Hypotheses

Based on prior work on phishing training efficacy [4], [5], [7], [10], [22], [23] and documented differences in risk perception [6], [9], we formulated four hypotheses to test in this work. Table 3 states each hypothesis, its basis, and its relation to prior work.

To summarize, **H1** represents the foundational assertion of the NIST Phish Scale framework [33]; **H2** and **H3** capture the intended effects of training interventions, with

TABLE 3: Summary of study hypotheses, their theoretical grounding, and relation to prior work. The hypotheses address the effects of phishing email difficulty, training, and training modality on detection and reporting behaviors.

Hypothesis	Theoretical and Empirical Basis	Relation to Prior Work
<b>H1 (Phish Scale Effect):</b> <i>Phishing emails with higher difficulty, as measured by the NIST Phish Scale, will exhibit higher click-through rates than those with lower difficulty.</i>	Grounded in the foundational NIST Phish Scale framework [33] and validation studies [39]. Supported by phishing cue research [34] and systematic email analysis [40].	The NIST Phish Scale has been validated in lab settings [39]. We provide the first large-scale enterprise validation.
<b>H2 (Training Impact — General):</b> <i>Training will increase the detection and reporting of phishing emails.</i>	Evidence of training efficacy, reduced click rates [22], improved detection accuracy [7], and short-term benefits [21]. Disputed by a recent study [4].	Tests whether training effects from small-scale studies [7], [22] replicate at scale, addressing the call for large-sample validation [10].
<b>H3 (Interactive Training Effect — Reporting):</b> <i>Interactive components (Trained+Exercise) will lead to significantly higher reporting rates than lecture-only training (Treatment A).</i>	From Franz <i>et al.</i> ’s modality taxonomy [10], and evidence of the effectiveness of active and feedback-rich training [8], [21], [22].	Builds on academic findings of interactive training benefits [21], [22]. We provide the first test at organizational scale.
<b>H4 (Interaction Effect):</b> <i>There will be a significant interaction between training modality and phishing difficulty, with interactive training being more beneficial for easier lures.</i>	Emerges from cognitive load and decision strategy research [19], difficulty-aware detection patterns [40], and findings on human limitations [6].	Addresses the gap identified by Franz <i>et al.</i> [10] regarding interaction effects between training and phishing complexity, largely unexplored empirically.

H3 specifically testing whether interactive modalities provide measurable benefits over passive learning approaches; and **H4** reflects our expectation that training effectiveness may diminish for highly sophisticated phishing attempts, as advanced phishing emails may exceed individuals’ detection capabilities regardless of training modality.

With respect to hypothesis novelty and prior testing, **H1** has been tested in small academic settings [39] but never at enterprise scale; **H2** and **H3** build on established training effectiveness research [7], [21], [22] but test these findings in a large-scale operational context using standardized difficulty measures; and **H4** is novel, addressing the research gap identified by Franz *et al.* [10] regarding interactions between training modality and phishing complexity, which has received minimal empirical attention in prior work.

### 3.2. Study Design

To test the hypotheses, we employed a between-subjects experimental design [35]. Subjects were randomly assigned to one of three experimental conditions using our organization’s training management system. Each condition reflected the phishing training they received (§3.3). Subjects were recruited (§3.4), trained, and later received simulated phishing emails (§3.5). We compared subjects’ performance in handling these phishing emails between the treatments using several standard metrics (§3.6).

### 3.3. Training Modalities

For the purposes of this research, the principal researcher collaborated with a cybersecurity training vendor who produces comprehensive security awareness training materials, phishing simulation exercises designed to augment those materials, and implementation assistance for organizational

deployment. This partnership enabled access to professionally developed training content and established phishing templates while maintaining the experimental rigor necessary for our evaluation.

The vendor offers three distinct training modalities, though we focused our evaluation on two of these approaches. The primary modality for cybersecurity awareness training consists of lecture-based videos that discuss the dangers of various social engineering attacks, including phishing recognition tactics, social engineering techniques, and organizational security protocols. Our vendor, along with many others in the industry, supplements this foundational training with interactive phishing exercises that present users with samples from a library of approximately 20 phishing templates. These exercises guide users through identifying suspicious elements and security deficiencies in emails that roughly correspond to the cues assessed by the NIST Phish Scale.

The vendor also provides a third modality referred to as “in-place training.” This approach occurs during phishing simulations when users click on malicious links, instead of being directed to a malicious website, they are immediately redirected to a training webpage that explains why the email they clicked should have been detected as fraudulent. However, decades of research into the effectiveness of embedded training environments [23] and recent comprehensive research by Ho *et al.* [4] convinced us that monitoring the results of in-place training was not necessary for our research objectives. Ho’s large-scale study across 19,500 healthcare employees found minimal effectiveness of embedded training approaches, providing thorough and convincing evidence that this modality offers limited practical value compared to the lecture-based and interactive approaches we chose to evaluate. We employed two distinct training approaches in collaboration with our cybersecurity training vendor, building on established tax-



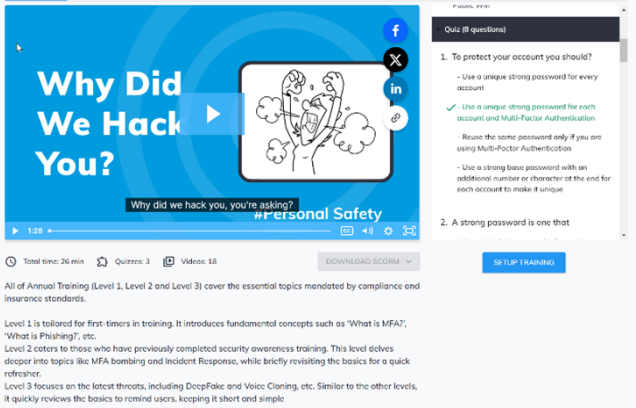


Figure 2: Traditional Training: Videos and Quiz

onomies of phishing interventions [10]. Both modalities fall within Franz *et al.*’s “training” category [10], specifically comparing passive lecture-based delivery Figure 2 against interactive approaches with embedded learning elements.

**Treatment A: Lecture-Based Training:** The lecture-style training consisted of fifteen brief instructional videos, each lasting 1-2 minutes, covering core cybersecurity awareness topics including phishing recognition tactics, social engineering techniques, and organizational security protocols [22]. The videos specifically addressed phishing identification strategies, emphasizing visual cues such as suspicious URLs, sender inconsistencies, and urgency tactics commonly employed by attackers. Following the video content, participants completed a knowledge assessment quiz reviewing key concepts from the training materials. This approach represents the traditional passive learning modality widely deployed in organizational cybersecurity training programs [23].

**Treatment B: Interactive Training (Training + Exercise):** The interactive group received the identical lecture-style training described above, plus an additional hands-on component consisting of simulated phishing exercises with immediate feedback. During these exercises, participants were presented with sample emails and asked to identify whether each message was legitimate or malicious. The system provided real-time feedback explaining the reasoning behind correct identification, highlighting specific phishing indicators, and reinforcing appropriate response behaviors such as using the “Report Phish” button in the email client [22]. This interactive component was designed to provide experiential learning opportunities that complement the passive video instruction, following research suggesting that active engagement enhances retention and behavioral change in cybersecurity training contexts [33].

**Control Group:** The control group received no cybersecurity training prior to the phishing simulation phase. This group was included to establish baseline susceptibility measurements. Control group participants received the same training materials after completion of the study to ensure organizational compliance with cybersecurity awareness re-

TABLE 4: Recruitment Pool and Completion Rates by Experimental Condition. Note the comparable drop-out rates by condition.

Condition	Recruited	Completed
Training Only (Treatment A)	6,758	6,023 (89.1%)
Training + Exercise (Treatment B)	6,764	6,026 (89.1%)
Control Group	529	462 (87.3%)
<b>Total</b>	<b>14,051</b>	<b>12,511 (89.0%)</b>

quirements. Organizational constraints compelled us to use a relatively small control group, as ACME Corp. required that the majority of employees receive some form of cybersecurity training for compliance purposes.

### 3.4. Subject Recruitment

**Study Context:** The subjects were employees at a U.S.-based international financial technology (“fintech”) firm. Participants represented a diverse cross-section of the organization, including employees from finance, technology, operations, customer service, and administrative functions. All participants were full-time employees with active email accounts who regularly use organizational communication systems and speak English.

We excluded participants who: (1) had inactive email accounts during the simulation period, (2) were on extended leave during the study period. Additionally, participants who received phishing simulations but had technical delivery failures (bounced emails, system errors) were excluded.

**Subject Recruitment:** This study leveraged organizational security training data collected as part of routine cybersecurity awareness training. The training was introduced through official organizational communications emphasizing its importance for maintaining security compliance and protecting organizational assets. This training was mandatory due to the organization’s need for regulatory compliance (cf. §2.1.2).

Data were collected from Dec. 2024 – Apr. 2025. Participants had a one month window to complete their assigned training modules, with automated reminders sent at regular intervals. Training completion was tracked and verified.

**Randomization:** Subjects were randomly assigned to each treatment. To mitigate bias by organizational role or function, we used stratified sampling [41] across departments to ensure representative distributions in each treatment group.

**Completion Rates:** Our data came from the 12,511 participants who completed the full experimental protocol. Table 4 shows the completion rates by treatment type. The comparable completion rates across conditions (ranging from 87.3% to 89.1%) suggest that our random sample was preserved between recruitment and completion.

**Ethical Considerations:** We acknowledge that employee consent in workplace security training cannot be considered

fully voluntary due to inherent power dynamics and employment dependencies [5]. However, several ethical safeguards were incorporated: all employees received security training regardless of group assignment, individual performance data was anonymized and never used for personnel decisions, and phishing simulation templates excluded content designed to cause psychological distress (§3.5).

The research protocol was reviewed by the academic institution’s IRB, which determined that analysis of operational security training data does not constitute human subjects research under federal guidelines, as it involves analysis of data collected for legitimate business purposes rather than research-specific data collection [5].

See Appendix 3 for an extended discussion of ethics.

### 3.5. Phishing Simulations

After training, all subjects received one simulated phishing email, of varying complexity. We first describe how we defined complexity using NIST Phish Scale, and then how we deployed the phishing simulation.

**3.5.1. Template Selection using NIST Phish Scale.** To establish controlled and consistent evaluation of participant responses, we systematically assessed and categorized the difficulty of each simulated phishing email using a structured expert review process informed by the NIST Phish Scale framework [33]. Phishing templates were initially selected from our cybersecurity training vendor’s library, which provided preliminary difficulty assessments (easy, medium, hard). We then independently evaluated these templates using the official NIST Phish Scale User Guide methodology to ensure accurate difficulty classification.

The assessment team comprised three members of ACME Corp.’s IT security staff, each with substantial experience in identifying real-world phishing attacks and analyzing user behavior in response to simulated campaigns. This multi-rater approach follows established practices for enhancing reliability and validity of stimulus assessment in cybersecurity research [19]. Two team members independently evaluated each phishing email template, rating the perceived human difficulty in correctly identifying the email as malicious using structured criteria derived from the NIST Phish Scale methodology.

Our implementation followed the official NIST Phish Scale User Guide [33], evaluating templates across two primary dimensions. For the *Phishing Cues* dimension, raters assessed observable indicators that might alert users to fraudulent content, including spelling errors, suspicious URLs, formatting inconsistencies, and sender authenticity markers [34]. Fewer or less obvious cues indicated higher difficulty ratings.

For the *Premise Alignment* dimension, raters evaluated how well each email’s subject matter, sender, and content aligned with typical organizational communications that ACME Corp. employees receive in their regular workflow. This assessment considered whether the email appeared to originate from familiar organizational systems

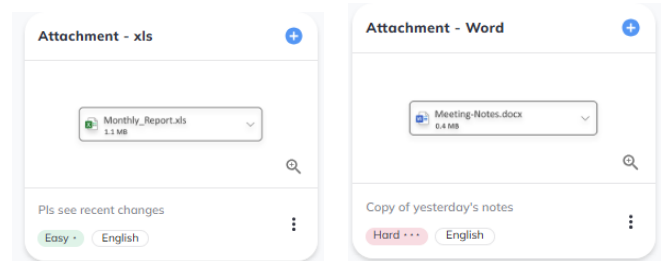
(e.g., HR portal, IT helpdesk, financial systems), addressed topics relevant to employee responsibilities, and used language and formatting consistent with legitimate internal communications [6]. Higher premise alignment—emails that closely resembled routine organizational correspondence—contributed to increased difficulty ratings.

Following independent assessment, the two raters compared their evaluations for each email template. Any disagreements regarding difficulty classification were discussed to identify the specific cues or contextual elements that led to differing conclusions. A third senior IT security expert then reviewed each email and the rationales provided by the independent raters, serving as a final validator to reach consensus difficulty ratings for each template. This collaborative validation process ensured that our phishing simulation represented a controlled range of difficulty levels essential for evaluating training effectiveness across standardized complexity measures [39].

**Difficulty Categorization and Distribution:** Following the NIST framework, we created a 3x3 classification matrix combining cue visibility (few-low, some-medium, many-high) with premise alignment (weak-low, medium, strong-high). Templates were assigned composite difficulty scores by summing both dimensions, with scores ranging from 2 (easy: high cues, low alignment) to 6 (hard: low cues, high alignment).

Our final simulation campaign included nine distinct phishing templates across three difficulty categories:

- **Easy:** 5,721 emails (45.7% of total simulations)
- **Medium:** 2,279 emails (18.2% of total simulations)
- **Hard:** 4,511 emails (36.1% of total simulations)



(a) “Easy”: low premise, high cues. (b) “Hard”: high premise, low cues.

Figure 3: Example templates. The “Easy” template has low premise alignment and high cues; the “Hard” is the inverse.

**3.5.2. Simulation Deployment.** Phishing simulations were distributed among participants within a three-month window after training completion. Each participant received one phishing email. Templates were randomly assigned to ensure balanced distribution across experimental conditions.

The deployed templates were constrained by organizational considerations. ACME Corp.’s management requested the exclusion of highly distressing content (e.g., emergency family situations, personnel security scenarios) to minimize employee psychological impact. This constraint reflects the

practical limitations of conducting research within operational environments, where employee welfare considerations influence experimental designs [9].

Technical delivery failures, including bounced emails and system errors, resulted in participant exclusion from analysis to maintain data integrity.

### 3.6. Metrics

Our outcome measures followed established practices in phishing research [22], [23].

- *Open Rate* captured the percentage of recipients who opened the simulated phishing email, measured via embedded tracking pixels, though we acknowledge this metric can be influenced by email client settings and security configurations [36].
- *Click-Through Rate (CTR)* measured the percentage of recipients who clicked on malicious links within the phishing emails, serving as our primary indicator of phishing susceptibility [7].
- *Reporting Rate* tracked the percentage of recipients who reported the suspicious email to the organization's security team using the integrated "Report Phish" button, representing proactive security behavior and successful threat recognition [9].

The simulation platform tracked these metrics as follows: email open rates were tracked via embedded tracking pixels, click-through rates were tracked via unique URLs, and reporting rates were tracked via integration with the organization's security incident reporting system.

### 3.7. Analysis

**Data Validation and Quality Assurance:** Prior to analysis, we conducted extensive data validation to ensure analytical integrity. This included verification of random assignment effectiveness across groups, assessment of technical delivery success rates, and examination of potential outliers or data anomalies. Participants with incomplete data due to technical failures, extended leave, or incomplete training were excluded from analysis to maintain dataset quality [39].

**Statistical Approach:** We employed two complementary analytical approaches to examine training effectiveness. Our primary analysis used two-way analysis of variance (ANOVA) to test main effects and interactions between training modality (Control, Training Only, Training + Exercise) and phishing difficulty level (Easy, Medium, Hard) on our outcome measures. Note that although ANOVA with binary outcomes violates standard normality assumptions, this approach aligns with established practices in phishing research [22], [23] and enables direct comparison with prior studies.

To address the limitations of ANOVA with binary data, we conducted confirmatory analyses using logistic regression models appropriate for binary outcomes (clicked/didn't click, reported/didn't report). The logistic regression models employed maximum likelihood estimation and provided

odds ratios for effect size interpretation. Both analytical approaches yielded consistent conclusions regarding training effectiveness, with ANOVA providing F-statistics for comparison with prior literature and logistic regression offering more statistically appropriate modeling of our binary outcomes. We report both approaches to provide comprehensive statistical evidence while acknowledging the methodological trade-offs inherent in each approach. All analyses were conducted using Python 3.9 with statsmodels and scikit-learn libraries for statistical modeling and cross-validation.

### 3.8. Threats to Validity

We encourage the reader to take several important limitations into consideration when interpreting the results of our study. We consider three threats to validity in our work: threats to construct, internal, and external validity [42].

**Construct Validity:** Construct validity is concerned with whether we operationalized the desired constructs of our study appropriately. In this study, we worked with two primary constructs: (1) the particular interpretations of the two different types of phishing training we deployed, which depend on our cybersecurity training vendor; and (2) our assessment of the difficulty of phishing lures.

It is possible that our cybersecurity vendor's training was simply lackluster, hence its failure to produce any significant improvement in subject performance. However, note that our results align with prior work that used different trainings. Meanwhile, to control for subjectivity in rating phishing lures, we applied the NIST Phish Scale and used multiple independent analysts to judge difficulty.

**Internal Validity:** Internal validity threats affect the reliability of conclusions we have about the relationship between the cause and the effect we observe. For our study, we used a stratified random sampling method in order to mitigate the potential biases related to job role and other organizational issues, allowing us to use that as a proxy for other factors such as technical background, educational attainment, *etc.* Due to the large sample sizes for our two different treatment groups, our study has significant statistical power; however, due to organizational constraints, we were required to limit the sample size for our control group.

**External Validity:** External threats impact the generalizability of our research. Our study was situated in the USA at a fintech company. The fintech industry context may limit generalizability to other sectors, as financial services employees may have different baseline levels of security awareness when compared to other industries. However, given the minimal effect sizes measured in our study, it is unclear how large of a threat this is. The study's focus on immediate post-training effects also limits our understanding of long-term training impact. Future research should examine whether the modest effects we observed persist over time or diminish as training recency decreases, or if they improve with volume or other dimensions.



TABLE 5: Summary of Training Effectiveness Results. Average rates show that trained groups performed similarly to control group, with modest differences across conditions.

Aspect	Results
Effect Size	Generally modest, consistent with recent large-scale findings [4]
<b>Click Rates</b>	
Overall Click Rate	10.4% across all participants
Control Group Average	9.8%
All Trained Average	10.5%
Training Only	10.6%
Training + Exercise	10.4%
<b>Report Rates</b>	
Overall Report Rate	9.6% across all participants
Control Group Average	8.9%
All Trained Average	9.7%
Training Only	9.5%
Training + Exercise	9.9%
<b>Statistical Effects</b>	
Training Effect on Click Rates	$p = 0.450$ (not significant)
Training Effect on Report Rates	$p = 0.417$ (not significant)
NIST Phish Scale Effect	$p < 0.001$ (highly significant)
Easy Difficulty	7.0% click rate
Medium Difficulty	8.7% click rate
Hard Difficulty	15.0% click rate
Interaction Effects	Revealed nuanced patterns with difficulty levels

## 4. Results

Our analysis examined the effectiveness of phishing training across 12,511 participants randomly assigned to three experimental conditions. Table 5 gives a summary of our key results. Next we present our findings organized by hypothesis, with detailed statistical analysis and effect size reporting for each. Our large-scale study revealed several key findings regarding the effectiveness of phishing training interventions and the predictive utility of the NIST Phish Scale.

### 4.1. H1: Phish Scale Effect

Training showed mixed effectiveness in our study. While we observed statistically significant effects for some outcomes, effect sizes were generally modest, consistent with recent large-scale findings suggesting limited practical impact of current training approaches [4]. Overall click rate

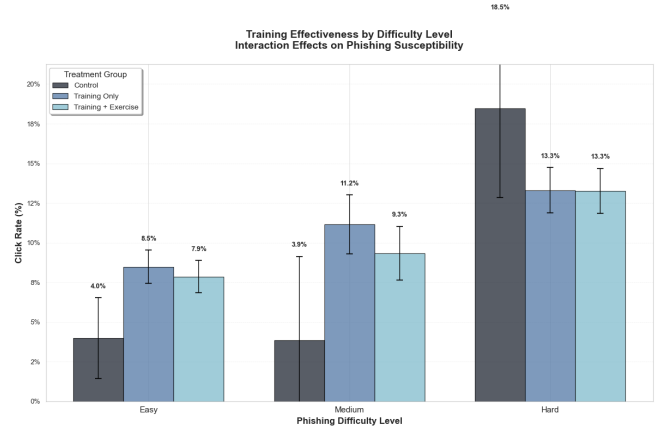


Figure 4: Comparison of effect sizes of phishing susceptibility. Bars show the variance explained by phishing difficulty (per the NIST Phish Scale) for each training intervention.

across all participants was 10.4% and overall report rate was 9.6%. Training on both click rates ( $p = 0.469$ ) and report rates ( $p = 0.448$ ) were not statistically significant in the primary analysis, though interactions with difficulty levels revealed more nuanced patterns.

The NIST Phish Scale successfully predicted user behavior, supporting H1. Phishing lure difficulty had a highly significant effect on click-through rates ( $F(2, 12086) = 41.415$ ,  $p < 0.001$ ,  $\eta^2 = 0.007$ ), validating the scale’s utility for measuring phishing complexity. Higher difficulty emails consistently produced higher click rates across all experimental conditions:

- Easy difficulty: 7.0% average click rate
- Medium difficulty: 8.7% average click rate
- Hard difficulty: 15.0% average click rate

This statistically significant difficulty gradient ( $F(2, 12086) = 41.415$ ,  $p < 0.001$ ,  $\eta^2 = 0.007$ ) confirms that the NIST Phish Scale effectively captures factors that influence user susceptibility to phishing attempts, providing empirical validation for its use in organizational contexts [33], [39].

### 4.2. H2: Overall Training Effects

Contrary to H2, training did not produce statistically significant improvements in either click reduction or reporting behavior. Main effect analysis showed no significant difference between trained and untrained groups for click rates ( $F(2, 12086) = 0.800$ ,  $p = 0.450$ ) or report rates ( $F(2, 12086) = 0.874$ ,  $p = 0.417$ ). This finding aligns with recent large-scale evidence suggesting minimal effectiveness of current training approaches [4].

However, examining treatment effects relative to control group baseline reveals some notable patterns:

#### Click Rate Changes (relative to control):

- Training Only: +4.5 percentage points for easy emails, +7.3 for medium, -5.2 for hard
- Training + Exercise: +3.9 percentage points for easy emails, +5.5 for medium, -5.2 for hard

#### Report Rate Changes (relative to control):

- Training Only: +0.1 percentage points for easy emails, +4.7 for medium, +1.9 for hard
- Training + Exercise: +0.8 percentage points for easy emails, +5.7 for medium, +1.7 for hard

The mixed direction of these effects suggests that training may have differential impacts across difficulty levels, though these did not reach statistical significance in the main effects analysis. The effect of training did not change significantly over time, as we observed similar patterns across the three-month simulation period.

### 4.3. H3: Interactive Training Effects

H3 was not supported. Interactive training (Training + Exercise) did not produce significantly higher reporting rates than traditional training alone. While the interactive group showed numerically higher reporting rates across difficulty levels, these differences were not statistically significant. The coefficient estimate for Training + Exercise relative to Training Only was modest and non-significant across all difficulty levels.

This finding contrasts with some prior work suggesting benefits of interactive training components [22], but aligns with recent large-scale studies showing limited training effectiveness regardless of modality [4].

### 4.4. H4: Training-Difficulty Interactions

H4 was not supported. We found no statistically significant interaction between training modality and lure difficulty for click rates ( $F(4, 12086) = 3.135, p = 0.014$ ) or report rates ( $F(4, 12086) = 0.517, p = 0.723$ ). While there was a marginal interaction effect for click rates, this primarily reflected the pattern that training appeared to increase clicks on easier emails while potentially reducing clicks on harder emails.

Coefficient analysis revealed some interesting interaction patterns:

- Training Only  $\times$  Hard: -9.67 percentage points ( $p = 0.002$ )
- Training + Exercise  $\times$  Hard: -9.09 percentage points ( $p = 0.004$ )

These significant negative interactions suggest that training may provide some protection against the most difficult phishing emails, even though overall training effects were not significant. It is worth considering that the training vendor created the phishing templates and training materials, which may introduce alignment between the cues emphasized in training and those present in the simulated emails. This alignment could artificially enhance the apparent effectiveness of training for certain email types, particularly those closely modeled after the vendor's examples. As a result, the observed reductions in click rates for hard emails may partially reflect this overlap rather than a generalizable improvement in user detection skills. Future studies should consider using independently developed phishing templates or varying template sources to better assess the robustness of training effects across diverse phishing scenarios.

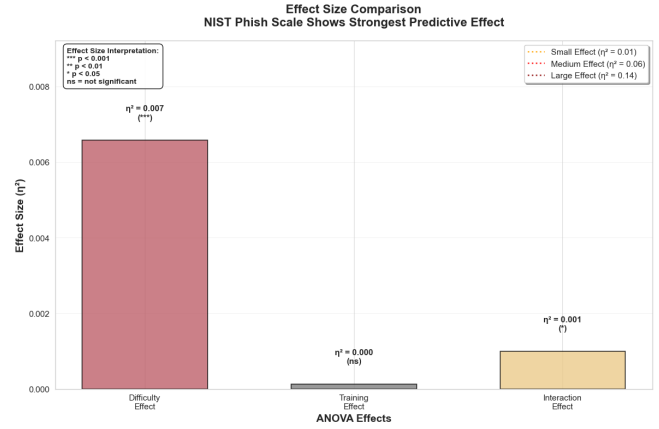


Figure 5: Comparison of effect sizes ( $\eta^2$ ) for key factors in phishing susceptibility. The NIST Phish Scale (difficulty) shows the strongest and only practically meaningful effect ( $\eta^2 = 0.007, p < 0.001$ ), while training and interaction effects are negligible ( $\eta^2 < 0.01$ ).

### 4.5. Effect Sizes and Practical Significance

While some individual coefficients reached statistical significance, effect sizes were generally small ( $\eta^2 < 0.01$  for all main effects), raising questions about practical significance [7]. The largest meaningful effect was the difficulty gradient ( $\eta^2 = 0.007$  for click rates), confirming the NIST Phish Scale's predictive validity.

Training effects, where present, were modest in magnitude. Even statistically significant interaction effects represented changes of less than 10 percentage points, suggesting that while training may provide some benefits, these are limited in scope and may not justify the substantial organizational investment required for comprehensive training programs [4].

These findings underscore the need for realistic expectations about training effectiveness and suggest that organizations should complement training with other defense mechanisms rather than relying on awareness programs as primary protection against phishing threats [10].

### 4.6. Other Measurements

The final question we considered was whether there was any relationship between training time and detection ability, *e.g.*, decay of training effectiveness. In our study, subjects received the phishing simulation email between 0–4 months after the training. Figure 6 shows the click rate by treatment group, discriminating by lure difficulty. We observe no discernible trend in the results as a function of time. Given the result of the visual inspection, we did not perform statistical tests of this aspect.

## 5. Discussion

Our large-scale field study provides empirical evidence for both the utility of the NIST Phish Scale and the limited

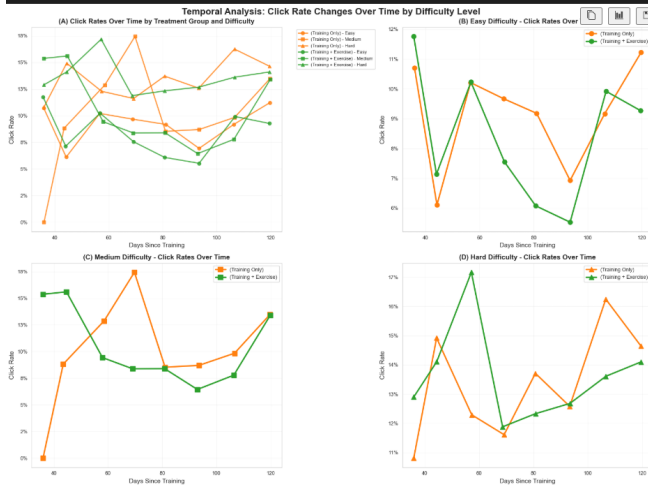


Figure 6: Temporal Impact on Detection. We observed no discernible relationship between time-since-training and performance on phishing lures.

effectiveness of current phishing training approaches. While we observed some statistically significant effects, the practical significance of these findings raises important questions about the return on investment for organizational security awareness programs.

### 5.1. Recommendations for Practice

**Confronting Training Effectiveness Realities:** Despite deploying well-structured training programs to over 12,000 employees, we found no statistically significant main effects of training on either click rates ( $p = 0.450$ ) or report rates ( $p = 0.417$ ). The small effect sizes we observed ( $\eta^2 < .01$  for all main effects) indicate that even statistically significant improvements translate to minimal operational impact. This aligns with recent large-scale evidence from Ho *et al.* [4] and challenges the optimistic assessments often promoted by cybersecurity training vendors.

The disconnect between vendor marketing claims and empirical evidence suggests that the substantial organizational resources typically invested in security awareness programs may not justify the modest improvements observed—even in best-case scenarios. Organizations must acknowledge these limitations to make informed investment decisions about cybersecurity resources.

**Layered Defense:** Our findings necessitate a fundamental shift in organizational phishing defense strategy. Rather than viewing training as a primary protection mechanism, organizations should treat security awareness programs as one component of a comprehensive layered defense strategy. The strong effect of phishing difficulty on user susceptibility ( $F(2, 12086) = 41.415, p < 0.001$ ) demonstrates that sophisticated attacks will likely succeed regardless of training efforts.

This evidence suggests several strategic reorientations: prioritizing passwordless authentication systems using

passkeys and FIDO-compliant hardware tokens like YubiKeys, especially for administrative roles; implementing comprehensive password manager deployment to reduce credential reuse vulnerabilities; enhancing email security controls specifically designed to catch sophisticated attacks; and developing incident response procedures that assume successful phishing will occur despite training efforts.

The reporting behavior improvements we observed, while modest, indicate that interactive training components may have some value in encouraging security-conscious behavior. However, organizations must weigh the practical significance of 1-2 percentage point improvements against the substantial additional resources required for interactive training delivery.

**Implementation Framework:** Organizations seeking to optimize their phishing defense should implement the following evidence-based framework:

**Technical Controls Priority.** Invest primarily in email filtering, attachment analysis, URL protection, and multi-factor authentication systems. Our findings suggest these provide more reliable protection than user training against sophisticated attacks.

**Targeted Training.** Rather than broad-based programs, implement focused training for specific roles, departments, or threat scenarios where evidence suggests training may be more effective. Use standardized difficulty assessment frameworks like the NIST Phish Scale for meaningful evaluation.

**Continuous Assessment.** Regularly evaluate the cost-effectiveness of training programs against alternative security investments, using standardized metrics that capture both behavioral outcomes and resource allocation efficiency.

This approach recognizes the sobering reality that current training approaches provide only modest protection against evolving phishing threats. The goal is not to eliminate training but to right-size its role within a comprehensive security strategy informed by empirical evidence rather than vendor marketing claims.

### 5.2. NIST Phish Scale Validation and Practical Utility

Our study provides strong empirical validation for the NIST Phish Scale’s predictive utility in organizational contexts. The highly significant effect of lure difficulty on click-through rates ( $F(2, 12086) = 41.415, p < 0.001$ ) demonstrates that the scale effectively captures factors influencing user susceptibility to phishing attempts. The clear gradient in click rates across difficulty levels—from 7.0% for easy emails to 15.0% for hard emails—confirms that the scale’s two-dimensional framework meaningfully differentiates phishing sophistication.

This validation extends the early work of Barrientos *et al.* [39] by demonstrating the scale’s effectiveness in a large operational environment. Unlike previous smaller-scale validations, our study shows that the NIST Phish Scale maintains its predictive power across diverse organizational

roles, departments, and user populations. This supports the scale's adoption as a standardized framework for evaluating phishing simulation programs and benchmarking organizational vulnerability assessments.

However, the scale's practical implementation revealed important considerations for practitioners. The process of systematically rating premise alignment across thousands of employees proved more complex than anticipated, requiring careful consideration of organizational communication patterns and role-specific email environments. Future implementations should account for these practical challenges when scaling the framework to enterprise environments.

### 5.3. Training Modality Effects and Interaction Patterns

While interactive training did not significantly outperform traditional lecture-based training in our primary analysis, the interaction patterns we observed provide nuanced insights into training effectiveness. The significant negative interactions between both training conditions and hard difficulty emails (Training Only  $\times$  Hard: -9.67 percentage points,  $p = 0.002$ ; Training + Exercise  $\times$  Hard: -9.09 percentage points,  $p = 0.004$ ) suggest that training may provide some protection against the most sophisticated phishing attempts. Conversely, the positive coefficients for training conditions with easier emails indicate that training may inadvertently increase vulnerability to simpler phishing attempts. This counter-intuitive finding aligns with concerns raised in prior work about potential negative effects of training programs [23]. One possible explanation is that training focuses employees' attention on detecting sophisticated attacks while reducing their vigilance toward obviously fraudulent emails.

These patterns underscore the complexity of human factors in cybersecurity training. Rather than uniformly improving detection across all threat types, training appears to redistribute user attention and response patterns in ways that may create new vulnerabilities. This highlights the importance of comprehensive evaluation frameworks that examine training effects across varied attack sophistication levels.

### 5.4. The Evolving Threat Landscape and Future Challenges

Our findings must be interpreted within the context of rapidly evolving phishing threats. The emergence of large language models and generative AI technologies is enabling attackers to create increasingly sophisticated and personalized phishing attempts that may render current training approaches obsolete. The traditional approach of teaching users to identify visual cues and suspicious content becomes less effective when attackers can generate near-perfect replicas of legitimate communications.

This technological evolution suggests that the cybersecurity community may need to fundamentally reconsider

human-centered defense strategies. Rather than expecting users to reliably detect sophisticated attacks, future approaches might focus on cryptographic verification systems, zero-trust architectures, in-person interaction (supporting return-to-office/RTO mandates) or other technical mechanisms that reduce reliance on human judgment.

The broader threat model is also expanding beyond traditional email phishing to include SMS phishing (smishing), voice phishing (vishing), and social engineering attacks using deep fakes across multiple communication channels. Training programs that focus narrowly on email-based threats may provide limited protection against this diversified attack landscape.

### 5.5. Implications for Theory and Practice

Our findings have significant implications for both cybersecurity theory and organizational practice. The limited training effectiveness we observed suggests that current human-centered defense paradigms may be fundamentally inadequate for addressing sophisticated phishing threats. The modest effect sizes ( $\eta^2 < 0.01$  for all main effects) indicate that even statistically significant improvements may not translate to meaningful operational risk reduction.

The interaction patterns we identified where training appeared to increase vulnerability to easier emails while potentially providing some protection against harder emails highlight the complex and sometimes counterintuitive effects of security awareness interventions. These findings suggest that training programs may redistribute rather than uniformly reduce user vulnerability, creating new attack vectors while addressing others.

For practitioners, our results support a realistic reassessment of training investments within broader security strategies. Organizations should view awareness programs as one component of layered defense rather than primary protection mechanisms. The substantial resources typically devoted to comprehensive training initiatives might achieve better risk reduction when reallocated toward technical controls, incident response capabilities, or more targeted interventions [10].

### 5.6. The Future of Human-Centered Cybersecurity

Our findings emerge at a particularly critical juncture in the evolution of cybersecurity. With the rapid advancement of generative AI enabled by large language models and other technologies, attackers are enabled to create increasingly sophisticated and personalized phishing attempts that may render traditional cue-based detection obsolete. When attackers can automate the creation of personalized phishing lures and generate near-perfect replicas of legitimate communications, it becomes unreasonable to expect users to reliably identify deceptive content.

The fundamental challenge lies in the inherent complexity of human factors in resistance to phishing attacks. Since these exploits target individual vulnerabilities, personal and contextual factors play an enormous role in individual user

susceptibility [19]. Even a single individual will vary dramatically throughout a day—users typically focus on content without paying attention to security, and are significantly impacted by cognitive constraints such as memory limitations and current cognitive load, all of which affect vigilance [43], [44]. Research has investigated variability in individual attentiveness based on factors like gender, age, and technical expertise, though the literature remains mixed [23], [38], [45]. Training represents only a single layer of defense against these multifaceted human vulnerabilities.

This technological evolution, combined with the complexity of human decision-making under cognitive load, suggests that the cybersecurity community must fundamentally reconsider where human-centered defenses lie in the multilayered approach organizations must take to cybersecurity. Future approaches may need to emphasize cryptographic verification systems, zero-trust architectures, and other technical mechanisms that can reduce reliance on the human firewall rather than attempting to enhance and upgrade the human firewall, which runs on hardware with a broad spectrum of variability.

The expanding threat landscape—encompassing SMS phishing, voice phishing, and multi-channel social engineering attacks—further challenges the viability of training-centric approaches focused primarily on email-based threats. As attack vectors diversify and sophistication increases, the cognitive burden placed on users to maintain vigilance across all communication channels may exceed reasonable human capabilities, particularly given the natural fluctuations in attention and decision-making quality that characterize normal human cognition. The role of human factors in resistance to phishing is complex.

## 6. Conclusion

This large-scale field study of 12,511 employees provides critical empirical evidence about the effectiveness of phishing awareness training in real-world organizational contexts. Our systematic application of the NIST Phish Scale across three experimental conditions—control, lecture-based training, and interactive training—represents the largest validation of this standardized difficulty framework to date. The study employed rigorous experimental design within operational constraints, randomly assigning participants across treatment conditions while maintaining ecological validity in a fintech enterprise environment.

Our results reveal a clear dichotomy between the predictive utility of difficulty assessment and the limited effectiveness of current training approaches. The NIST Phish Scale demonstrated strong predictive power, with phishing difficulty significantly affecting user susceptibility ( $F(2, 12086) = 41.415$ ,  $p < 0.001$ ), showing click rates ranging from 7.0% for easy emails to 15.0% for hard emails. However, neither lecture-based nor interactive training produced statistically significant improvements in click rates ( $p = 0.450$ ) or reporting behavior ( $p = 0.417$ ), with effect sizes below 0.01 for all main training effects.

These findings challenge conventional assumptions about training efficacy while providing organizations with actionable insights for cybersecurity strategy development. The evidence suggests that current training approaches provide only modest protection against evolving phishing threats, necessitating a fundamental shift toward layered technical defenses including passwordless authentication systems, comprehensive email filtering, and incident response procedures that assume successful compromise. Rather than viewing training as a primary defense mechanism, organizations should integrate awareness programs as supplementary components within comprehensive security architectures that emphasize technical controls over human judgment for critical security decisions.

## References

- [1] Z. Alkhalil, C. Hewage, L. Nawaf, and I. Khan, “Phishing Attacks: A Recent Comprehensive Study and a New Anatomy,” *Frontiers in Computer Science*, vol. 3, p. 563060, Mar. 2021, <https://www.frontiersin.org/articles/10.3389/fcomp.2021.563060/full>.
- [2] F. L. Greitzer, W. Li, K. B. Laskey, J. Lee, and J. Purl, “Experimental Investigation of Technical and Human Factors Related to Phishing Susceptibility,” *ACM Transactions on Social Computing*, vol. 4, no. 2, pp. 1–48, Jun. 2021, <https://dl.acm.org/doi/10.1145/3461672>.
- [3] F. Menges, R. Paccagnella, G. Pellegrino, and M. Johns, “Contrasting and synergizing CISOs and Apprentices: A persona-based analysis of Organizational Security Cultures,” in *Proceedings of the CHI Conference on Human Factors in Computing Systems*. ACM, May 2024, pp. 1–18.
- [4] G. Ho, A. Mirian, E. Luo, K. Tong, E. Lee, L. Liu, C. A. Longhurst, C. Dameff, S. Savage, and G. M. Voelker, “Understanding the Efficacy of Phishing Training in Practice,” *IEEE Symposium on Security and Privacy*, 2025.
- [5] J. M. Haney, S. M. Furman, and Y. Acar, “Compliance to impact: Tracing security awareness programme evolution,” *Computers & Security*, vol. 128, p. 103123, 2023.
- [6] H. Abroshan, J. Devos, G. Poels, and E. Laermans, “Phishing Happens Beyond Technology: The Effects of Human Behaviors and Demographics on Each Step of a Phishing Process,” *IEEE Access*, vol. 9, pp. 44 928–44 949, 2021, <https://ieeexplore.ieee.org/document/9380285/>.
- [7] A. Sumner, X. Yuan, M. Anwar, and M. McBride, “Examining Factors Impacting the Effectiveness of Anti-Phishing Trainings,” *Journal of Computer Information Systems*, vol. 62, no. 5, pp. 975–997, Sep. 2022.
- [8] L. Schöni, V. Carles, M. Strohmeier, P. Mayer, and V. Zimmermann, “You Know What? - Evaluation of a Personalised Phishing Training Based on Users’ Phishing Knowledge and Detection Skills,” in *Proceedings of the 2024 European Symposium on Usable Security*. ACM, pp. 1–14. [Online]. Available: <https://dl.acm.org/doi/10.1145/3688459.3688460>
- [9] A.-K. Doing, E. Bárbaro, F. Van Der Roest, P. Van Gelder, Y. Zhauniarovich, and S. Parkin, “An analysis of phishing reporting activity in a bank,” in *Proceedings of the 2024 European Symposium on Usable Security*. Karlstad Sweden: ACM, Sep. 2024, pp. 44–57. [Online]. Available: <https://dl.acm.org/doi/10.1145/3688459.3688481>
- [10] A. Franz, V. Zimmermann, G. Albrecht, K. Hartwig, C. Reuter, A. Benlian, and J. Vogt, “SoK: Still Plenty of Phish in the Sea — A Taxonomy of User-Oriented Phishing Interventions and Avenues for Future Research,” *Proceedings of the Seventeenth Symposium on Usable Privacy and Security*, pp. 1–22, 2021.



- [11] A. Ferreira and S. Teles, "Persuasion: How phishing emails can influence users and bypass security measures," *International Journal of Human-Computer Studies*, vol. 125, pp. 19–31, May 2019, <https://linkinghub.elsevier.com/retrieve/pii/S1071581918306827>.
- [12] A. Das, S. Baki, A. El Aassal, R. Verma, and A. Dunbar, "SoK: A Comprehensive Reexamination of Phishing Research From the Security Perspective," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 1, pp. 671–708, 2020, <https://ieeexplore.ieee.org/document/8924660/>.
- [13] G. Desolda, L. S. Ferro, A. Marrella, T. Catarci, and M. F. Costabile, "Human Factors in Phishing Attacks: A Systematic Literature Review," *ACM Computing Surveys*, vol. 54, no. 8, pp. 1–35, Nov. 2022, <https://dl.acm.org/doi/10.1145/3469886>.
- [14] A. Aleroud and L. Zhou, "Phishing environments, techniques, and countermeasures: A survey," *Computers & Security*, vol. 68, pp. 160–196, 2017.
- [15] A. E. Aassal and R. M. Verma, "Spears Against Shields: Are Defenders Winning the Phishing War?" in *Proceedings of the 4th ACM Workshop on Security Information Workers*, 2019, pp. 3–10, <https://dx.doi.org/10.1145/3309182.3309191>.
- [16] T. Bikkur, M. Nikitha, A. Vajja, K. Harshitha, and J. Rani, "Optimized Machine Learning Algorithm to classify Phishing Websites," in *2022 IEEE International Conference on Electronics, Computing and Communication Technologies (CONECCT)*, 2022, <https://dx.doi.org/10.1109/ICEARSS53579.2022.9752223>.
- [17] P. K. P. M. M. A. Butavicius, M., "Breaching the human firewall: Social engineering in phishing and spear-phishing emails," in *ACIS 2015 Proceedings*, 2015, pp. 1–10.
- [18] A. Burns, M. E. Johnson, and D. D. Caputo, "Spear phishing in a barrel: Insights from a targeted phishing campaign," *Journal of Organizational Computing and Electronic Commerce*, vol. 29, no. 1, pp. 24–39, 2019.
- [19] J. S. Downs, M. B. Holbrook, and L. F. Cranor, "Decision strategies and susceptibility to phishing," in *Proceedings of the Second Symposium on Usable Privacy and Security - SOUPS '06*. Pittsburgh, Pennsylvania: ACM Press, 2006, p. 79, <http://portal.acm.org/citation.cfm?doid=1143120.1143131>.
- [20] C. I. Canfield, B. Fischhoff, and A. Davis, "Better watch out: Comparing phishing metacognition with real emails," *Learning and Metacognition*, vol. 14, no. 3, pp. 343–362, 2019, <https://doi.org/10.1007/s11409-019-09197-5>.
- [21] N. Marshall, D. Sturman, and J. C. Auton, "Exploring the evidence for email phishing training: A scoping review," *Computers & Security*, vol. 139, p. 103695, Apr. 2024, <https://linkinghub.elsevier.com/retrieve/pii/S0167404823006053>.
- [22] A. Carella, M. Kotsoev, and T. M. Truta, "Impact of security awareness training on phishing click-through rates," in *2017 IEEE International Conference on Big Data (Big Data)*. IEEE, 2017, pp. 4458–4466.
- [23] D. D. Caputo, S. L. Pfleeger, J. D. Freeman, and M. E. Johnson, "Going Spear Phishing: Exploring Embedded Training and Awareness," *IEEE Security & Privacy*, vol. 12, no. 1, pp. 28–38, Jan. 2014, <http://ieeexplore.ieee.org/document/6585241/>.
- [24] P. Bitrián, I. Buil, S. Catalán, and D. Merli, "Gamification in workforce training: Improving employees' self-efficacy and information security and data protection behaviours," *Journal of Business Research*, vol. 179, p. 114685, Jun. 2024, <https://linkinghub.elsevier.com/retrieve/pii/S0148296324001899>.
- [25] S. Hart, A. Margheri, F. Paci, and V. Sassone, "Riskio: A Serious Game for Cyber Security Awareness and Education," *Computers & Security*, vol. 95, p. 101827, Aug. 2020, <https://linkinghub.elsevier.com/retrieve/pii/S0167404820301012>.
- [26] M. Franklin, T. Folke, and K. Ruggeri, "Optimising nudges and boosts for financial decisions under uncertainty," *Palgrave Communications*, vol. 5, no. 1, p. 113, Oct. 2019, <https://www.nature.com/articles/s41599-019-0321-y>.
- [27] U. D. of Health and H. Services, "Security standards: Technical safeguards - hipaa security rule," 2023, <https://www.hhs.gov/hipaa/for-professionals/security/laws-regulations/index.html>.
- [28] E. Parliament and C. of the European Union, "General data protection regulation (gdpr) - article 39: Tasks of the data protection officer," 2016, <https://gdpr-info.eu/art-39-gdpr/>.
- [29] P. C. I. S. S. Council, "Payment card industry data security standard (pci dss) v4.0," 2023, [https://www.pcisecuritystandards.org/document\\_library](https://www.pcisecuritystandards.org/document_library).
- [30] I. O. for Standardization (ISO), "Iso/iec 27001:2022 - information security, cybersecurity and privacy protection," 2022, <https://www.iso.org/standard/27001>.
- [31] N. I. of Standards and T. (NIST), "Framework for improving critical infrastructure cybersecurity, version 2.0," 2024, <https://www.nist.gov/cyberframework>.
- [32] M. Steves, K. Greene, and M. Theofanos, "Categorizing human phishing difficulty: A Phish Scale," *Journal of Cybersecurity*, vol. 6, no. 1, p. tyaa009, Jan. 2020, <https://academic.oup.com/cybersecurity/article/doi/10.1093/cybsec/tyaa009/5905453>.
- [33] S. Dawkins and J. Jacobs, "NIST Phish Scale user guide," Gaithersburg, MD, p. NIST TN 2276, Nov. 2023. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/TechnicalNotes/NIST.TN.2276.pdf>
- [34] M. Canham, S. Dawkins, and J. Jacobs, "Not All Victims Are Created Equal: Investigating Differential Phishing Susceptibility," in *Augmented Cognition*, D. D. Schmorow and C. M. Fidopiastis, Eds. Cham: Springer Nature Switzerland, 2024, vol. 14694, pp. 3–21, [https://link.springer.com/10.1007/978-3-031-61569-6\\_1](https://link.springer.com/10.1007/978-3-031-61569-6_1).
- [35] D. Hillman, Y. Harel, and E. Toch, "Evaluating organizational phishing awareness training on an enterprise scale," vol. 132, p. 103364. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0167404823002742>
- [36] D. Antoniou, "A systematic method to execute simulated phishing tests," *Journal of Phishing Prevention*, vol. 4, no. 2, pp. 112–118, 2015, [https://scholar.google.com/scholar?hl=en&as\\_sdt=0%2C5&q=A+systematic+method+to+execute+simulated+phishing+tests&btnG=](https://scholar.google.com/scholar?hl=en&as_sdt=0%2C5&q=A+systematic+method+to+execute+simulated+phishing+tests&btnG=).
- [37] M. Canham, C. Posey, and M. Constantino, "Phish Derby: Shoring the Human Shield Through Gamified Phishing Attacks," *Frontiers in Education*, vol. 6, p. 807277, Jan. 2022, <https://www.frontiersin.org/articles/10.3389/feduc.2021.807277/full>.
- [38] M. De Bona and F. Paci, "A real world study on employees' susceptibility to phishing attacks," in *Proceedings of the 15th International Conference on Availability, Reliability and Security*. Virtual Event Ireland: ACM, Aug. 2020, pp. 1–10, <https://dl.acm.org/doi/10.1145/3407023.3409179>.
- [39] F. Barrientos, J. Jacobs, and S. Dawkins, "Scaling the Phish: Advancing the NIST Phish Scale," in *HCI International 2021 - Posters*, C. Stephanidis, M. Antona, and S. Ntoa, Eds. Springer International Publishing, vol. 1420, pp. 383–390. [Online]. Available: [https://link.springer.com/10.1007/978-3-030-78642-7\\_52](https://link.springer.com/10.1007/978-3-030-78642-7_52)
- [40] R. Wright, S. Johnson, and B. Kitchens, "Phishing Susceptibility in Context: A Multilevel Information Processing Perspective on Deception Detection," *MIS Quarterly*, vol. 47, no. 2, pp. 251–270, 2023, <https://dx.doi.org/10.25300/misq/2022/16625>.
- [41] T. Sutter, A. S. Bozkir, B. Gehring, and P. Berlich, "Avoiding the Hook: Influential Factors of Phishing Awareness Training on Click-Rates and a Data-Driven Approach to Predict Email Difficulty Perception," *IEEE Access*, vol. 10, pp. 100 540–100 565, 2022, <https://ieeexplore.ieee.org/document/9893815/>.
- [42] R. Verdecchia, E. Engström, P. Lago, P. Runeson, and Q. Song, "Threats to validity in software engineering research: A critical reflection," *Information and Software Technology*, vol. 164, p. 107329, 2023.

- [43] CK. Fallon, JA. Baweja, JY. Yun, ND. Thompson, ZH. Shaw, and DL. Arendt, "Phishing in the Wild: An Ecologically Valid Study of the Phishing Tactics and."
- [44] N. C. Ebnner, D. M. Ellies, T. Linn, H. A. Rochaa, H. Yang, S. Dommaraju, A. Solliman, D. L. Wooddard, G. R. Turner, R. N. Spreng, and D. S. Olliveira, "Risk of sensitivity to online disappointment in old age," *Journals of Gerontology - Series B Psychological Sciences and Social Sciences*, vol. 75, no. 3, pp. 522–533, 2020, <https://doi.org/10.1093/geronb/gby036>.
- [45] D. Conway, R. Taib, M. Harris, K. Yu, S. Berkovsky, and F. Chen, "A qualitative investigation of bank employee experiences of information security and phishing," in *Thirteenth Symposium on Usable Privacy and Security (SOUPS 2017)*, 2017, pp. 115–129.

## Outline of Appendices

The appendix contains the following material:

- Appendix 2: Detailed information about phishing lures.
- Appendix 3: Ethical considerations in study design and execution.

### 1. Acknowledgments

We thank the employees of ACME Corp. for their participation in this study.

### 2. Lures

See Table 6 for our analysis of the phishing lures used in this study.

### 3. Ethical Considerations

The research protocol was reviewed by the IRB of the academic institution, which determined that analysis of operational security training data does not constitute human subjects research under federal guidelines, as it involves analysis of data collected for legitimate business purposes rather than research-specific data collection [5].

ACME Corp. employees are informed during onboarding that cybersecurity training and periodic phishing assessments are mandatory conditions of employment, consistent with industry standards for financial services organizations [9]. All participants had previously acknowledged organizational policies regarding security training participation through standard employment agreements. This approach aligns with established practices in organizational cybersecurity research, where security training represents a legitimate business function rather than experimental manipulation [22].

We acknowledge that employee consent in workplace security training cannot be considered fully voluntary due to inherent power dynamics and employment dependencies [5]. However, this limitation applies broadly to organizational security research and reflects the practical constraints of conducting ecologically valid studies in operational environments. Recent work by Haney *et al.* demonstrates that meaningful organizational security research requires balancing employee autonomy concerns with the legitimate need to evaluate security program effectiveness [5].

The study design incorporated several ethical safeguards to minimize potential harm to participants. First, all employees received security training regardless of group assignment, with control group participants receiving identical training materials after the simulation phase to ensure organizational compliance requirements were met. Second, individual performance data was anonymized for analysis and never used for personnel decisions or disciplinary actions, consistent with research ethics guidelines for workplace studies [9]. Third, phishing simulation templates were constrained by organizational welfare considerations, excluding content designed to cause psychological distress (*e.g.*,

emergency scenarios, financial hardship themes) following established ethical frameworks for security awareness research [23].

We recognize that employees who underperform in security training may face additional training requirements as part of standard organizational security protocols. However, these consequences exist independently of our research activities and represent standard business practices for maintaining organizational security posture. Our analysis focused on aggregate training effectiveness patterns rather than individual performance evaluation, ensuring that research participation did not create additional risks beyond those inherent in routine security training programs [10].

This approach reflects broader challenges in conducting ethical security research within operational constraints, where the alternative of laboratory-based studies may lack ecological validity necessary for informing real-world security practices [7]. Recent scholarship emphasizes that organizational security research must balance competing ethical considerations while advancing knowledge that can improve organizational and societal cybersecurity resilience [5].

TABLE 6: NIST Phish Scale Assessment of Phishing Templates. Premise Alignment scores range from 1 (Low) to 3 (High). Lure Scores range from 1-3 based on phishing cue assessment. Cue Count represents the total number of suspicious indicators identified in each template. Difficulty categories were determined by combined assessment of phishing cues and premise alignment following NIST Phish Scale methodology.

Template Name	Difficulty	Cue Count	Lure Score	Premise Score	Premise Alignment Explanation
Attachment - Word	Hard	14	2	2	Moderate alignment: Document sharing is plausible, but placeholder fields reduce realism.
Vulnerabilities — Office	Hard	18	3	1	Low alignment: Scenario is vague and not role-specific; unlikely context.
Attachment - XLS	Easy	14	2	2	Moderate alignment: XLS financial data plausible, but lacks sender/context.
Loom - Recording Shared	Hard	14	2	2	Moderate alignment: Loom recordings common, but template fields weaken trust.
Device Non-compliant - Microsoft	Medium	15	3	3	High alignment: Matches expected MDM behavior and terminology.
Fax via Voicemail Office	Easy	13	2	1	Low alignment: Fax delivery uncommon in modern workplace.
Email Not Delivered — Office	Easy	13	2	3	High alignment: Bounce messages are routine and widely expected.
Email Not Delivered — Google	Easy	16	3	3	High alignment: Gmail bounce notices with technical detail are believable.
From Marketing: Collab Request	Easy	14	2	3	High alignment: Internal collaboration messages from marketing are common.
Accounting Team Amex Charges — Office	Easy	16	3	3	High alignment: Finance inquiries during month-end are familiar.
HSA Use Balance!	Easy	17	3	2	Moderate alignment: HSA reminders are expected, but lacks personalization or sender.
Mobile Device Restricted — Outlook	Medium	16	3	3	High alignment: Device quarantine messages match common security flows.
Tax Workshop	Easy	16	3	2	Moderate alignment: HR tax webinars plausible but weakened by generic formatting.
Mentor Program	Hard	15	3	3	High alignment: Internal mentorship invitations are believable and well-framed.
HR - Performance Review — Outlook	Hard	16	3	3	High alignment: Outlook calendar invites for reviews are routine.
Changes to Healthcare Form	Hard	15	3	3	High alignment: Healthcare form updates are common in HR communications.
Device Non-compliant - Google	Medium	16	3	3	High alignment: Technical restriction notices align with enterprise environments.
Vulnerabilities — Google	Hard	16	3	1	Low alignment: Scenario lacks specificity and role-relevant context.