# "An incident may have resulted in a suspected data compromise": Impact of Data Breach Notification Terminology

Xiaoxin Shen, Ann Zhang, Xinyi Hu, Yixuan Wang, Mingjie Chen, Kai Sze Luk

*Carnegie Mellon University*

## Abstract

We conducted an exploratory study of user reaction to common terminology used in Massachusetts data breach notifications. We surveyed 99 Prolific participants and asked them to evaluate segments of breach notifications. We found that overall, participants had high security awareness and had a solid understanding of the concepts of phishing and what might be considered as personally identifiable information. We also identified a lack of relevant details in the data breach notification as a potential barrier that stops users from taking remediation actions after a security breach. Our findings suggest that while users have the common impression that data breaches occur often, they are not yet familiar with relevant data breach notification laws and possible remediation steps.

## 1. Introduction

With the increase in frequency of reported data breaches across different states in the US, there is a need to examine how companies are communicating this information with impacted users [1]. Data breach notifications can be defined as "letters written by organizations that have experienced a breach to alert users that sensitive personal information has been accessed, lost, or stolen, which usually puts users at a higher-than-normal risk of identity theft." [2]. Currently, all 50 states in the United States require organizations and government entities to notify users when personally identifiable information has been exposed or organizations' data security measures have been breached [3].

Aside from being a legal requirement, data breach notifications also impact user behavior following the event. Previous research showed that upon receiving a data breach notification, customer spending decreased significantly, with some customers migrating to use alternative services that were not breached [4]. Researchers also found that while customers generally had high levels of information security awareness, they were relatively unaware of related legislative regulations and were only partially aware of possible remediation strategies [5][6].

In this work, we created a data breach notification based on the data breach notification archive from the state of Massachusetts, and evaluated its readability using Flesch Reading Ease Score (FRES) and Flesch-Kincaid Grade Level (FGL) analysis, before deploying an online survey where we asked 99 participants throughout the United States to evaluate the notification [7]. We found that while participants had high awareness of information security risks, their knowledge of data breach notification laws is lacking, and a lack of relevant details in a breach notification can act as a barrier for users to take remedial action after the data breach event.

## 2. Methodology

### 2.1. Experiment Design

Using postal mail data breach notifications from banks in February 2022 in Massachusetts, we developed a list of high frequency words using Optical Character Recognition (OCR) and NLP analysis [8]. We decided to use data breach notification letters from Massachusetts due to its publicly available breach notification archive. While California has a similar archive, the state law only requires breached entities to submit a notice to the Attorney General's office if impacted California residents exceed 500, making it a much smaller archive compared to Massachusetts, where breached entities are required to submit a notice regardless of number of individuals impacted [9][10]. From this analysis, the high frequency terms are "report, credit, secure, freeze, inform, card, request, identity, place, and account."

Inspired by the high frequency terms, we selected relevant segments from different data breach statements and combined those into one coherent notification (See Appendix A - Survey Protocol Part 1.1 - 1.2). Participants were asked to imagine themselves as clients of UMC Bank, who had just received a data breach notification. After reviewing the statements, participants were asked to respond to open-ended questions about how likely they thought their information was leaked, who might now have access to their information, what they thought a "phishing event" is, who is at fault in this data breach and how well they felt their data was protected by UMC bank. Then, participants saw another statement from UMC vouching to protect their information better, and were asked to outline what is considered "personally identifiable protected data," what remediation steps they would take and how likely they think there would be a future data breach at UMC. We again asked if

participants felt their data was protected to see if written commitment towards data protection would impact user impression of the bank. Then, we asked about participants' actual experience with breach notifications, and if applicable, what steps they took after receiving such communication. This allowed us to compare what participants claim they would do to their actual behavior after receiving a data breach notification. The survey concludes with demographic and technological background questions and was designed to take 10 minutes to complete. Our protocol was approved by Carnegie Mellon University's Internal Review Board. We did not collect any personally identifiable information from the participants.

## 2.2. Participant Recruitment

We limited participant age to 18 and older and those who are fluent in English as the original statements were written in English. As all 50 states in the United States had some form of data breach notification law and the statement was developed from the Massachusetts' data breach notification archive, we decided to limit participation to those located in the United States. Since our survey required extensive reading, we limited survey response to those who are accessing Prolific via desktop computers. We recruited 101 participants from Prolific, with an average completion time of 10 minutes and 4 seconds, compensating each participant $1.59 USD for a complete submission.

We piloted our study protocol with 22 participants prior to launching it on Prolific. Our original study protocol focused on standalone sentences used in data breach notifications, but as the sentences did not provide sufficient background information, many participants focused on the lack of background information in their responses. The UMC bank scenario was developed to help pivot the participants' focus back onto the notification terminology.

## 2.3. Data analysis.

To analyze the readability of the data breach notification we developed for our survey, we used Flesch Reading Ease Score (FRES) and the Flesch Grade Level (FGL) as proposed by Zou et al [7]. FRES and FGL show to what extent a breach notification is comprehensible by the general public. Zou et al. also used the Gunning Fog index (FOG), a grade-level-based metric that factors in complex words (those containing three or more syllables).

Our qualitative analysis includes responses from 99 participants, with 2 responses being removed as they were suspected bot-replies. The two responses had identical typos and the replies did not answer our questions meaningfully. We utilized emergent coding to identify common themes and sentiment and we double-coded every question to increase reliability. Due to time constraints, we were unable to perform Cohen's Kappa test for reliability.

# 3. Results

## 3.1. Readability and Structure of our notification

To assist the reader with recognizing the dangers posed by the breach and what measures to take, a data breach notice should use clear language structured in an accessible format. For each notice used in our survey, we evaluated its readability, expected reading duration, and the usage of structural headers.

We identified several readability issues through FRES and FGL analysis. FRES evaluates texts on a 0 - 100 scale, with higher scores indicating easier readability. The statement used in Survey Protocol Part 1.1(S1) has a FRE Score of 53.0, with the second statement (S2) having a score of 50.0, meaning both text was considered Fairly difficult (50-59) to read (See Appendix A for full statement).

Converting the FRES to FGL, S1 received a FGL score of 11.2, and S2 being 12.3. This indicates that to be able to understand half of the statement, an 11th to 12th grader's reading aptitude is needed. Prior literacy research suggests that materials addressed to the general public should aim for a junior-high reading level, with a FGL between 7 to 9 [7]. According to our demographic data, 7.9% of our participants have a high school degree, with 92.1% of our participants having a high school degree or higher. We expect that while the statements may be challenging, our participants have sufficient reading aptitude to understand these statements.

## 3.2. User understanding of technological terms

In order to evaluate participants' knowledge of technical terms, we asked participants to define "What is a phishing event?" then compared user generated replies to our definition of phishing, being "a digital social engineering technique where the perpetrator masquerades as another entity and attempts to trick users into revealing sensitive information." [11] We identified four critical aspects of phishing, being Social Engineer (manipulation, social engineer, trick), Masquerade (attacker pretends to be another entity, impersonation), Information (Revealing personal information), Format (form of solicitation, email, SMS, faulty link) and checked if participants mentioned any of the four aspects outlined.

Overall, we found that 6.1% of participants did not know what phishing is, and on average participants mentioned 1.8 out of the four aspects outlined, with 7.1% mentioning all four. 70.9% of participants mentioned revealing sensitive information, with 48.5% also referring to the impersonation aspect involved in phishing. Overall, participants showed a high awareness of information security with varying degrees of accuracy in understanding what phishing is.

We also asked participants to outline what they considered as "personally identifiable protected information" (PII) and compared user generated responses to forms of PII listed by government agencies[12][13][3]. Mostly, participants' understanding of PII aligned with that identified in the law. It is interesting to note that while Email is mentioned by CCPA and DHS as a form of PII, only 5.1% of participants agreed as such.

### 3.3. User behavior after receiving notification.

After reading our notification, 29.3% of participants said they would contact the bank, with 53.5% saying they would replace their cards and 19.2% would consider switching banks if they "didn't feel like the situation is in control". 19.1% of participants would monitor their credit and account transaction history closely, 20.2% would change log-in credentials, with 5.1% asking the bank for advice on this situation. Notably, only 2% of participants said they would "probably [do] nothing."

In Section 2 of our survey, we asked if participants had any actual experience with data breach notifications, and if applicable, what actions they took after receiving one. Similarly, around 20% reported changing their credentials and 9% monitored their accounts. However, the percentage of participants who reported taking no action soared to 19.2%. Reasoning varied across participants, with one participant accrediting their inaction to being "already enrolled in identity protection," while another complained that the notification "did not say which card … was compromised, making it very difficult … to do anything about it."

### 3.4. User impression of data breach event post notification

After reading S1, the top reason that participants cited that suggests they were personally impacted was that the bank reached out. One participant explained that "The bank sent me a letter specifically listing all of the info that may have been compromised. It's bad PR for them so the breach must have been really bad." This statement also highlights another common reason that gave the impression data was leaked, being that the notification provided an extensive list of what information may have been compromised.

For those who were unsure if they were impacted, top reasons included that the statement used "vague language" and seems to be "mass sent." 13.1% of participants singled out "may have happened" and "some of our customers" as reasons to suggest the email might have been a "cursory notice." A few participants also questioned the validity of our scenario, cushioning answers with "assuming that the message is legitimate … and not a phishing attempt," or suggesting that our notification could be "a bank scam."

Participants were also asked about who they thought now had access to their information, with 75.8% of the responses

pointing towards hackers, fraudsters and anyone who was related to the incident. 12.1% of participants also mentioned their information was probably available to "anyone who is willing to buy it," indicating data brokers. When asked who was at fault in this breach, 62.6% of the participants thought the bank was responsible for the breach, with 17.2% mentioning the attacker. Participants commented that both "the perpetrator and the institution who was hacked" were at fault. 14.1% mentioned that the IT security team was also partially responsible for not "being protected enough."

### 3.5. User confidence in organization post notification

At the end of S1, we asked participants how confident they felt that their data was protected and why they felt that way. Then, we showed participants S2, which was a written statement from UMC dedicated towards data protection and asked them to re-evaluate their confidence. For the initial question, 70.7% of participants reported they felt their data was not protected, with 8.1% of participants reporting feeling protected. The remaining 21.2% was somewhere in the middle, feeling that their data was somewhat protected.

The most common reason participants cited for feeling not protected was that the leak happened. 34.3% of participants believed that the fact that the leak was possible in the first place indicated that the bank was not protecting their data. 9.1% of participants cited the attack type as the reason for feeling not protected, because phishing was considered a "low-level attack" that the bank should be capable of defending against. For participants who felt their data was protected, the most common reason was trust in the bank or organizations in general. 6.1% of participants stated that they trusted banks to protect their data, while 4.0% cited a trust in organizations and companies in general to do their best for security, not necessarily only banks.

Interestingly, participants cited attacks being common both as a reason that their data was protected and not protected. 11.1% felt that as attacks are common, there was no way their data was protected, whereas 4.0% believed that the leak would occur despite the bank protecting their data because attacks are common.

After seeing S2, 40.4% of participants kept the same confidence rating, while 23.2% reported decreased confidence and 36.4% increased confidence. 12.1% of participants cited S2 as new information that made them feel better about their data being protected, while 15.2% of participants reported S2 as a reason they felt their data was less protected than they previously believed. Those who were convinced by the additional statement explained that after a breach, they believed that the bank would improve their security and the probability of another breach in the future was low. For those not so convinced, they felt that the bank either didn't give enough information about security

measures, were not taking enough additional security measures, or there was nothing the bank could do to make the situation better.

We also asked participants to rate their confidence that a data breach would not occur again in the future, and explain why. The vast majority of participants (85.9%) reported low confidence that a data breach would not occur again. Out of these participants, 24.2% of participants felt that once a data breach occurred, it was likely to happen again, and 20.2% mentioned that due to the frequency of data breaches, it was likely that one would happen again. The main reason participants cited for having confidence that a breach would not be repeated was that the bank would make security improvements after this breach.

## 4. Discussion

### 4.1. High Information Security Awareness of Phishing

Our study results confirm findings from previous research that users have a high information security awareness [5]. Our participants at minimum were aware that phishing was an attempt to gain personal information from them, with almost half of the participants knowing that phishing utilized some sort of impersonation technique. Our sample scenario and notification was also questioned as possible phishing attempts by some participants, showing high vigilance in security awareness.

### 4.2. Discrepancy between Security Awareness and Knowledge of Data Breach Notification Laws

The results also highlight the gap between consumer security awareness and their actual knowledge of data breach notification laws, which confirms findings from previous research [5]. The top reason for suspecting that user information was compromised was that the bank reached out, with participants either directly saying or implying institutions would avoid doing so unless the situation is dire. Breached entities who compromised personally identifiable information are required by law to contact their consumers in all 50 states across the country, yet there seems to be a lack of user expectation in terms of being informed by breached entities should a data breach occur.

### 4.3. Deviation between User Claims and Actual User Behavior

While many users said that they would consider changing service providers after receiving a breach notification, their actual behavior deviated from their claims. Building on previous research where consumers were only partially aware of possible remediation strategies [6], we found that users had legitimate barriers that deterred them from taking remediation action. The lack of detailed information from the breach notification makes it difficult for users to isolate the main account that has been compromised.

### 4.4. Limitations

The statements we used in this survey had a FRES of 51.5 and a FGL of 11.75, which is more difficult to read compared to the recommended FGL score between 7 - 9 [7]. Our sample breach notification database was also limited to the data breach notification archive from Massachusetts, in the future it would be interesting to explore and compare different data breach notifications between states, and compare the difference in information mandated by state government.

Another limitation is that while we attempted to pivot participants' focus onto terminology, as we only provided segments of a standard data notification, some participants still asked for more details, which a typical data breach notification would have included. Another factor that we did not consider was that in our scenario we did not provide a send-date. The survey was deployed on April 18th 2022, with many participants naturally assuming that this was the date when they received the notification. The breach incident in our scenario was dated to have taken place between August 24th 2021 - October 14th 2021. This gives the impression that users were notified around 6 months after the discovery of the incident, which may have led to some bias in our qualitative data, as some participants saw the long period between the time of breach (August 24th 2021 - October 14th 2021) and the date of notification (April 18th 2022) as a sign of untrustworthiness or incompetence from the bank.

### 4.5. Future Work

An observation that could inspire future research would be the low number of participants who considered email addresses as a form of PII. The Department of Homeland Security, Massachusetts State Laws and California Consumer Privacy Act all mention email as a form of PII, so it would be prudent to explore email and whether it is considered as a form of personal or public information from a user perspective.

Users citing attacks being common as both reasons that their data is protected and not protected is perhaps an observation into current user impression on cybersecurity. More research is needed in this area to explore the relationship between this belief and possible user inaction after receiving a data breach notification, as this may encourage the user to take more action to protect themselves, or none at all.

## Acknowledgements

## References

[1] K. Grindal, "What Works? Measuring the Efficacy of Cyber Policy Interventions with Quasi-Experiments," Social Science Research Network, Rochester, NY, SSRN Scholarly Paper ID 3893092, Jul. 2021. doi: 10.2139/ssrn.3893092.

[2] J. R. Veltsos, "An Analysis of Data Breach Notifications as Negative News," *Business Communication Quarterly*, vol. 75, no. 2, pp. 192–207, Jun. 2012, doi: 10.1177/1080569912443081.

[3] National Conference of State Legislatures, "Security Breach Notification Laws," 2022. https://www.ncsl.org/research/telecommunications-and-information-technology/security-breach-notification-laws.aspx

[4] R. Janakiraman, J. H. Lim, and R. Rishika, "The effect of a data breach announcement on customer behavior: Evidence from a multichannel retailer," *Journal of Marketing*, vol. 82, no. 2, pp. 85–105, 2018, doi: 10.1509/jm.16.0124.

[5] J. Nield, J. Scanlan, and E. Roehrer, "Exploring Consumer Information-Security Awareness and Preparedness of Data-Breach Events," *Library Trends*, vol. 68, no. 4, pp. 611–635, Spring 2020, doi: http://dx.doi.org/10.1353/lib.2020.0014.

[6] Z. Hassanzadeh, R. Biddle, and S. Marsen, "User Perception of Data Breaches," *IEEE Transactions on Professional Communication*, vol. 64, no. 4, pp. 374–389, Dec. 2021, doi: 10.1109/TPC.2021.3110545.

[7] Y. Zou, S. Danino, K. Sun, and F. Schaub, "You `Might' Be Affected: An Empirical Analysis of Readability and Usability Issues in Data Breach Notifications," in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, New York, NY, USA, May 2019, pp. 1–14. doi: 10.1145/3290605.3300424.

[8] Office of Consumer Affairs and Business Regulation Massachusetts, "Data Breach Notification Letters | Mass.gov," 2022. https://www.mass.gov/archive/data-breach-notification-letters (accessed May 09, 2022).

[9] General Court of the Commonwealth of Massachusetts, *General Law - Part I, Title XV, Chapter 93H, Section 1*. 2007. Accessed: May 08, 2022. [Online]. Available: https://malegislature.gov/Laws/GeneralLaws/PartI/TitleXV/Chapter93H/Section1

[10] State of California Department of Justice, "Search Data Security Breaches," *State of California - Department of Justice - Office of the Attorney General*, 2022. https://oag.ca.gov/privacy/databreach/list (accessed May 08, 2022).

[11] NIST, "phishing - Glossary | CSRC," 2022. https://csrc.nist.gov/glossary/term/phishing (accessed May 07, 2022).

[12] Department of Homeland Security, "What is Personally Identifiable Information? | Homeland Security," 2021. https://www.dhs.gov/privacy-training/what-personally-identifiable-information (accessed May 07, 2022).

[13] State of California Department of Justice, "California Consumer Privacy Act (CCPA)," *State of California - Department of Justice - Office of the Attorney General*, Oct. 15, 2018. https://oag.ca.gov/privacy/ccpa (accessed Feb. 13, 2022).

## Appendix

A - Survey Protocol

Part 1.1 - Statement 1

Scenario:

Imagine that you are a client of UMC bank, with whom you have multiple accounts, a credit card and debit card. You have just received a letter from the bank and when you opened the document, it begins with the following:

"This notice is to inform you of a suspected data compromise. We have reasons to believe that some of our customers may have had their data compromised.

What happened: An incident occurred between 8/24/2021-10/14/2021 may have resulted in the disclosure of your information due to a bank vendor phishing event.

What information was involved: According to our records, the information involved in this incident was related to your loan and may have included your first and last name, address, account number, credit/debit account number and routing number."

With this in mind, please answer the following:

1. How likely do you believe it is that your information was leaked? (1 = Not Likely, 10 = Very Likely)

2. Looking at your prior answer, what made you select that rating?

3. What information do you think was disclosed?

4. Who do you think now has access to your information?

5. What is a "phishing event"?

6. Who is at fault in this data breach?

7. How confident are you that your data is protected? (1 = Not Confident, 10 = Fully Confident)

8. Looking at your answer above, what factors contributed (or did not contribute) towards your confidence?

Part 1.2 - Statement 2

Imagine that the letter concludes with the following:

"UMC Bank takes its obligation to safeguard personally identifiable protected data entrusted to us very seriously and therefore deem it necessary to bring this situation to your attention. We want to inform you of what we are doing to protect you and what you can do to protect yourself.

You may visit a branch for a new card, or you may request we mail your new UMC debit card in about 10-20 business days."

Follow-up questions:

1. What information do you believe is considered "personally identifiable protected data"?

2. What steps would you take after receiving this letter?

3. How confident are you that a data breach will not occur again in the future? (1 = Not Confident, 10 = Fully Confident)

4. Looking at your answer above, what factors contributed (or did not contribute) towards your confidence?

5. How confident are you that your data is protected? (1 = Not Confident, 10 = Fully Confident)

6. Looking at your answer above, what factors contributed (or did not contribute) towards your confidence?

Part 2 Follow up

1. Have you ever received communication (letter, email) that included similar statements?
   a. Yes
   b. No
   c. Unclear

2. When was the last time you received a notice containing similar statements?
   a. In the last 6 months
   b. In the last 12 months
   c. Not Applicable
   d. Prefer not to answer
   e. Other (Please specify)

3. What (if anything) did you do after receiving communication that your information has been compromised?

Part 3 Demographics

1. What is your gender?
   a. Male
   b. Female
   c. Prefer not to answer
   d. Other (Please specify)

2. What is your age group?
   a. 18-25 years old
   b. 26-35 years old
   c. 36-45 years old
   d. 46-55 years old
   e. 55+ years old
   f. Prefer not to answer

3. What is your race/ethnicity?
   a. Hispanic or Latino - A person of Cuban, Mexican, Puerto Rican, South or Central American, or other Spanish culture or origin regardless of race.
   b. White (Not Hispanic or Latino) - A person having origins in any of the original peoples of Europe, the Middle East, or North Africa.
   c. Black or African American (Not Hispanic or Latino) - A person having origins in any of the black racial groups of Africa.
   d. Native Hawaiian or Other Pacific Islander (Not Hispanic or Latino) - A person having origins in any of the peoples of Hawaii, Guam, Samoa, or other Pacific Islands.
   e. Asian (Not Hispanic or Latino) - A person having origins in any of the original peoples of the Far East, Southeast Asia, or the Indian Subcontinent, including, for example, Cambodia, China, India, Japan, Korea, Malaysia, Pakistan, the Philippine Islands, Thailand, and Vietnam.
   f. American Indian or Alaska Native (Not Hispanic or Latino) - A person having origins in any of the original peoples of North and South America (including

Central America), and who maintain tribal affiliation or community attachment.

    g. Two or More Races (Not Hispanic or Latino) - All persons who identify with more than one of the above five races.

    h. Prefer not to answer

4. What is your occupation?

    a. CS/IT-related

    b. Somewhat related to IT

    c. Other (Please specify)

    d. Prefer not to answer

5. Which of the following best describes your highest achieved education level?

    a. Some High School

    b. High School Graduate

    c. Some college, no degree

    d. Associates degree

    e. Bachelor's degree

    f. Graduate degree (Masters, Doctorate, etc.)

    g. Other (Please specify)

    h. Prefer not to answer

Part 4 - Technology

1. Please select all of the digital devices that you own:

    a. Smart Phone

    b. Smart Watch

    c. Laptop Computer

    d. Desktop Computer

    e. Smart Home devices (Alexa, Amazon Echo and similar devices)

    f. Gaming Device

    g. VR Headset

    h. iPad and/or electronic tablets

    i. Other (Please specify)

2. Please rate this statement: I can solve most of my own technical problems

(1 = Strongly Disagree, 10 = Strongly Agree)

3. Please rate this statement: I know about a wide range of different technology

(1 = Strongly Disagree, 10 = Strongly Agree)

4. Please rate this statement: I keep up with technological news

(1 = Strongly Disagree, 10 = Strongly Agree)