






## Regular Article

## Can you spot a scam? Measuring and improving scam identification ability

Elif Kubilay<sup>a</sup>  Eva Raiber<sup>b</sup>  Lisa Spantig<sup>c,a,\*</sup>  Jana Cahlíková<sup>d</sup>  Lucy Kaaria<sup>e,f</sup><sup>a</sup> University of Essex, United Kingdom<sup>b</sup> Aix-Marseille Université, CNRS, AMSE, France<sup>c</sup> RWTH Aachen University, Germany<sup>d</sup> University of Bonn, Germany<sup>e</sup> University of Nairobi, Kenya<sup>f</sup> HOPAWI, Kenya

## ARTICLE INFO

Dataset link: <https://doi.org/10.7910/DVN/HUYYQZ>

## JEL classification:

Codes

D14

D18

G53

O12

## Keywords:

Consumer protection

Consumer fraud

Digital financial services

Scam susceptibility

Scam education

Kenya

## ABSTRACT

The expansion of digital financial services leads to severe consumer protection issues such as fraud and scams. As these potentially decrease trust in digital services, especially in developing countries, avoiding victimization has become an important policy objective. In an online experiment, we first investigate how well individuals in Kenya identify phone scams using a novel measure of scam identification ability. We then test the effectiveness of scam education, a commonly used approach by organizations for fraud prevention. We find that common tips on how to spot scams do not significantly improve individuals' scam identification ability, i.e., the distinction between scams and genuine messages. This null effect is driven by an increase in correctly identified scams and a decrease in correctly identified genuine messages, indicating overcaution. Additionally, we find suggestive evidence that genuine messages with scam-like features are misclassified more often, highlighting the importance of a careful design of official communication.

## 1. Introduction

The expansion of digital financial services (DFS) has increased access to financial services, both in developed and developing countries (e.g., Pazarbasioglu et al., 2020; Balyuk, 2022). With this increase in DFS, consumer protection issues are also on the rise (Garz et al., 2021). One major issue is fraud. Fraud is detrimental to consumers both in terms of direct monetary costs and indirect costs such as erosion of trust in financial services (Guiso et al., 2008; Gurun et al., 2017; Johnson et al., 2019), loss of confidence in financial matters (Brenner et al., 2020), and mental health problems including depression and stress (DeLiema et al., 2020; Financial Institution Regulatory Authority, 2015). One common type of fraud is phone scams using text messages or calls. The goal of scammers is to trick consumers into sending money or revealing private information such that their accounts can be accessed. Since scammers often target random phone numbers, all segments of society who have a phone and use basic DFS are at risk.

The negative effects of committed and attempted fraud might go beyond the scope of financial services. If there is a high level of

mistrust due to the prevalence of scams, individuals might tend to ignore messages, undermining the effectiveness of digital messages as a communication tool. This can have important implications for the functioning of markets, the provision of information, and public service delivery. For example, SMS-based communication has been used to reduce frictions in rural labor and agricultural markets (Fabregas et al., 2019). Messages have also been used to enhance individuals' knowledge and health behaviors (Holst et al., 2021; He et al., 2023), and to motivate bureaucrats (Dustan et al., 2023). These examples all rely on employers, employees, citizens, and bureaucrats to open, read, and consume the content of messages. In contexts where mobile phones are the only way to reach large shares of the population, fear of fraud may hinder communication with these groups.

The existing recipe for avoiding consumers' scam victimization is to pursue education and awareness campaigns. Yet, do educational campaigns indeed improve people's ability to detect scams and do they influence how genuine messages from e.g. banks or telecommunication providers are perceived? An important obstacle to evaluating the effect

\* Corresponding author at: University of Essex, United Kingdom.

E-mail address: [lisa.spantig@essex.ac.uk](mailto:lisa.spantig@essex.ac.uk) (L. Spantig).

of education campaigns is quantifying the relevant outcome metrics. Consumers under- or misreport fraud attempts and victimization: They might not be able to recognize all types of fraud, differentiate genuine offers from scams, or remember all instances of fraud attempts (Chen et al., 2018). Moreover, victims often feel shame and guilt and do not report scams to avoid potential stigma (Burke et al., 2022). Therefore, we argue that a policy-relevant metric is the ability to identify fraud attempts and confidence in this ability. Even if only a few individuals are direct victims of fraud, the inability to recognize fraud or the lack of confidence in this ability may impede market participation.

In this paper, we study susceptibility to scams and the effectiveness of a light-touch scam education in Kenya. First, we develop a novel measure for an individual's scam identification ability (SIA) and confidence in their ability. For this, we collect actual scams and official communication that circulate in Kenya. Second, we test experimentally if common tips for scam detection improve SIA and confidence. We focus on Kenya, Africa's leader in digital infrastructure and mobile money use (Koyama et al., 2021). At the same time, the country suffers from increasing rates of phone scams, which by now represent the most often cited consumer protection issue (Blackmon et al., 2021).

In an online survey (N = 1000) we show respondents 12 different messages and ask them to indicate whether these messages are scam or not. Each classification decision is followed by a confidence rating. The messages include both common scams and genuine messages sent by, e.g., banks or telecommunication companies in Kenya. After having classified the first six messages, a random half of the respondents receive tips on fraud prevention that are commonly provided by banks or telecommunication companies. These tips warn consumers about "scam markers", which include (i) typos and grammar mistakes, (ii) an unknown sender, (iii) a shortened link, and (iv) requests for private information such as pin codes or passwords. Ideally, these tips help respondents become better at distinguishing scams from genuine messages. However, it is also possible that tips about scams make respondents more cautious and hence more likely to classify any given message as scam. The latter would make it harder for service providers to communicate with their clients.

We find that on average tips do not increase scam identification ability. This null effect arises because while respondents in the treatment group are more likely to correctly identify scams, they are also less likely to correctly identify genuine messages. On average, tips appear to make consumers more cautious, i.e., more likely to classify any given message as scam. Moreover, receiving tips makes respondents significantly more confident in their classification decisions. The increase in confidence could be concerning as average SIA does not increase in our study and overconfidence has been found to be correlated with victimization (McAlvanah et al., 2015). However, we find suggestive evidence that higher confidence is associated with better SIA at the individual level.

Looking deeper, we find a more nuanced result depending on whether a given message contains a "scam marker". First, tips increase the number of correctly identified scams, irrespective of whether a scam marker is present in the message. This suggests that tips indeed make people more cautious. Second, a scam marker in a genuine message increases the likelihood of this message being classified as a scam. Part of the null effect of tips thus seems to be driven by official messages that look like scams. This highlights that tips need to be specific enough to unambiguously increase SIA and that non-scam communication should avoid features that are commonly cautioned against in educational campaigns.

To test whether the results change when money is at stake, a random half of our sample receives incentives for each correct classification. Results show that incentives do not lead to better SIA, and there is no interaction effect with tips. While we find some indications that those who receive both tips and incentives exert more effort, these individuals are not performing better. This implies that our measure can be used as an unincentivized survey measure. We also investigate treatment

effect heterogeneity and illustrate an important shortcoming of such education interventions. Tips appear to be effective in increasing SIA only for more experienced DFS users and those with higher education. Less educated participants do not benefit from tips, suggesting that it is difficult to design universally-helpful communication.

This study relates to several strands of literature. First, we contribute to the nascent literature on financial fraud in developing countries. For example, Ensminger and Leder-Luis (2022) and Andersen et al. (2022) study the detection of fraud in foreign aid, whereas Garz et al. (2021) summarize the consumer protection challenges of the expansion of DFS. Different types of fraud have been documented in various settings: fraudulent smartphone apps in India (Fu and Mishra, 2022), phone scams as the most prominent consumer protection issue in Kenya (Blackmon et al., 2021), and agent misconduct in Ghana (Annan, 2022a,b). Here, we focus on phone scams, develop a measure of SIA based on actual scams and official messages, and show that information makes individuals more careful on average but does not increase their SIA.

Second, we contribute to the literature studying the causal effects of educational interventions on fraud susceptibility. This literature has mostly focused on phishing attacks (e.g., Sheng et al., 2007) but also studied telemarketing schemes (Scheibe et al., 2014), and investment scams (Burke et al., 2022). In general, tips and information may decrease fraud susceptibility, especially among better-educated individuals (Burke et al., 2022). Our study suggests that previous findings, albeit focusing on different types of scams, appear to hold for phone scams and in a developing country setting. Additionally, we make a methodological contribution by making mistakes costly for half of our sample and show that this does not alter results. Participants' intrinsic motivation to correctly classify messages appears to be high enough.

Third, we contribute to a large literature documenting correlates of fraud susceptibility and victimization (see Moustafa et al., 2021; Norris et al., 2019, for recent reviews). This literature studies samples from developed, Western countries. The most common demographic characteristics that have been found to matter are gender and age. Additionally, financial knowledge (Engels et al., 2020), as well as risk aversion, curiosity, and the level of trust (Chen et al., 2018) are associated with fraud susceptibility. We find similar results in the Kenyan context. Women and less experienced DFS users have lower SIA. Additionally, we show that these groups do not differentially benefit from tips. These results imply that unless information provision is more targeted at specific groups of the population, such policy interventions are unlikely to close existing gaps in SIA.

## 2. Background

Kenya is a leading market for digital financial products and services (Koyama et al., 2021). Often, solutions are tested in Kenya and then rolled out to other countries in the region. With near-universal phone penetration and use of DFS in Kenya, almost all adults are at risk of phone scams. In a representative survey of active DFS users, 56% reported they had been contacted by scammers in the past six months, most commonly by phone (Blackmon et al., 2021). Scam reports are prevalent among all demographic groups. Interestingly, 68% of users with tertiary education reported scam attempts, compared to only 50% among those with at most secondary education. This suggests that less educated consumers might not recognize all scam attempts and/or might be less willing to report them.

Given the high prevalence of scams, it is not surprising that 90% of the adult population is concerned about fraud when using digital services (Koyama et al., 2021). In terms of direct costs of victimization, recent numbers show a positive correlation between the depth of digital services use and the amount lost due to fraud.<sup>1</sup> This implies that with

<sup>1</sup> Koyama et al. (2021) find that over the past three years, more advanced users lost more than twice as much as the basic digital services users due to fraud.

increased use of DFS, more people will be at risk of suffering from unexpected losses that might be difficult for them to absorb. Regarding indirect costs, 71% of the self-employed report limiting their usage of DFS due to concerns about fraud (Koyama et al., 2021), indicating loss of trust.

Qualitative interviews and scam examples from social media that we collected show that scammers try to trick individuals into transferring money or to obtain personal information to either access accounts or steal the identity of the victim. Scammers often impersonate bank and telecommunication agents, relatives or friends. A variety of different scams exist, from fake loan or investment offers to prizes for which money has to be sent upfront to take advantage of these “opportunities”. “Erroneous contact” is another common scam in which the sender pretends to have sent money or sensitive information and either asks for the money to be transferred back or for the enticing information to be ignored. In the latter cases, the primary goal is to start a conversation for more sophisticated social engineering.

While Kenya has passed digital safety policies and laws, and has established the office of the Data Protection Commissioner, the problem of fraud cannot be solved by regulation alone. Technological innovations such as biometric identification can help protect identities and accounts, but the human factor also needs to be addressed. In other contexts, it appears that financial knowledge is associated with lower susceptibility to fraud (Engels et al., 2020), which might explain the general popularity of educating consumers to raise awareness of and resilience to fraud (DeLiema et al., 2020; Engels et al., 2020).

### 3. Measuring scam identification ability

To build our measure of scam identification ability (SIA), we obtained information about ongoing scams from different sources. First, using a social media analytic tool, we collected public posts from Twitter sent between January 2020 and June 2021 from a Kenyan location. We kept the posts that were sent from an individual account and related to phone scams based on topic clustering. We further restricted the sample to contain screenshots of text messages which, after removing duplicates, left us with 116 tweets. Additionally, we conducted a survey in the largest Kenyan fraud-detection Facebook group in September 2021. Members of the group were asked to submit examples of both scam and official messages and calls. Participants submitted 922 examples, of which about 62% were scams. As the type of messages, i.e., scam or official, is self-reported and might be subject to error, we hired two research associates to independently classify 516 messages (including 116 from Twitter) and assert their confidence. In cases where two coders' classification did not match, a third research associate was asked to make a classification.

We focus on SMS scams and non-scam text messages to use examples verbatim.<sup>2</sup> To generate variation in our SIA measure, we construct a database of “ambiguous” messages, where either the two coders did not agree or the average confidence rating of the classification was low. All ambiguous messages were discussed within the research team and with experts if needed.<sup>3</sup> From this set of ambiguous messages, we randomly select 13 scams and seven official messages, stratified by topic. We turn these messages into vignettes by equalizing the visual appearance and pilot the 20 vignettes in two small convenience samples (N = 39). We select the 12 final vignettes based on the classification decisions and confidence of pilot participants.

<sup>2</sup> Recalled protocols of calls were incomplete, similar to examples of SMS that were not copy-pasted or submitted as a screenshot. Administering the vignettes in a written context (in our online survey) allows us to keep the mode of perception close to real life.

<sup>3</sup> We describe the process of building the measure in more detail in Appendix C.

Our measure consists of two blocks with four scam vignettes and two official messages each.<sup>4</sup> The blocks are presented in random order, and the messages are randomized within each block. We refer to the block shown first as “block 1” and the one shown second as “block 2”. For each block, we measure SIA as the share of correctly classified messages. We also examine separately whether individuals classify scams and non-scams correctly. As we are more interested in the former, we decided to include more scam than non-scam messages in our measure. For each vignette, participants indicate whether this is a scam or not (binary choice). Afterwards, a scale appears on the same page and asks participants to rate their confidence in their classification on a five-point Likert scale where the higher values indicate higher confidence.

### 4. Experimental setting

We measure SIA in an online survey in which we also administer an education treatment to estimate the causal effect of scam tips on the ability to distinguish fraudulent from genuine messages.

#### *Tips treatment*

Educational campaigns aim at raising awareness and providing tips on how to distinguish scam and non-scam communication (e.g., “Safaricom will only SMS you from MPESA and Safaricom”) or on how to behave (e.g., “never share your PIN”). These campaigns are often run visually on billboards or social media. Therefore, to capture available information on fraud prevention, we collected examples of tips using Twitter and qualitative data. We condense the five most common pieces of information into one infographic (see Fig. 1). To avoid information overload and ensure that all tips are read, we animate the graphic, such that the participants see one bullet point at a time. Participants go through this animation at their own speed. On average, they spent 1.12 min (SD = 0.67) reviewing the tips.

We randomize scam education at the individual level and provide it to 50% of our sample. We administer the treatment between the two blocks of vignettes, which allows us to assess individuals' SIA level prior to tips treatment. It is important to note that, as in real life, we do not distinguish between information being new or serving as a reminder.

#### *Incentive treatment*

In contrast to real life where mistakes can be costly, our participants may exert less effort. We hence cross-randomize a robustness treatment in which we pay 10 KES for each correctly classified message. Half of our sample receives incentives in both blocks. Different from the tips treatment, incentives may thus influence all classifications. We opted to pay incentives from the beginning such that participants who receive both tips and incentives can focus on understanding the main treatment between the two blocks. Finally, the incentive treatment allows us to explore whether using incentives is essential to elicit scam identification ability.

<sup>4</sup> For non-scams, we focus on official communication by banks, Safaricom as the provider of MPESA, and other telecom providers. As we exclude circumstantial clues from our design, personal messages from family and friends cannot be unambiguously classified as non-scam. As an unknown sender is the most obvious clue for a scam, we vary whether the sender is shown in the vignette. See Table A1 for an overview of all vignettes and Figure A1 for a visual example.

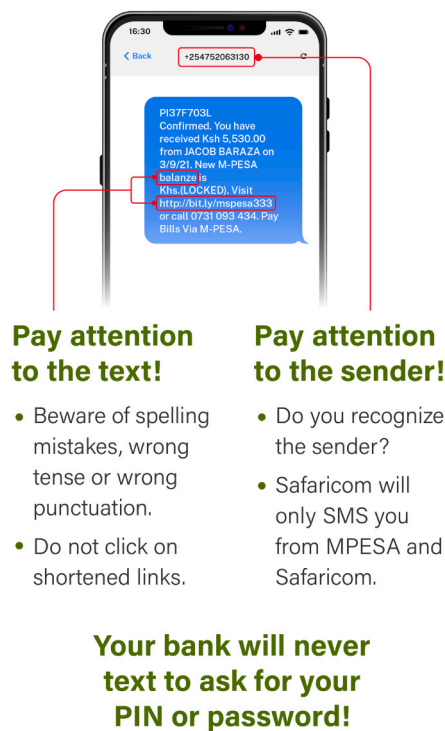


Fig. 1. Tips treatment.

Notes: Tips treatment was designed based on commonly communicated tips in Kenya. The graphic was “animated”, such that the pieces of information would be shown step-by-step. Participants clicked through this animation at their own speed, i.e., they hit the “continue” button five times before they see the overall graphic.

### Online survey and the sequence of events

After written consent and questions on demographics, phone ownership and usage, participants are shown a definition of scams and told that their task is to identify scam messages. They do not receive information about the number of vignettes or the fraction of fraudulent messages. Before starting the first block, participants in the incentive treatment learn about the payment for correct classification. After the first block, participants in the tips treatment go through the animated infographic. Nobody receives feedback on their SIA measured in the first block. Afterwards, everyone proceeds with the second block, followed by questions regarding the use of DFS, scam experiences, and an attention check. At the end of the survey, participants learn the number of correctly identified messages and those in the incentive treatment also see the corresponding bonus payment.

### Procedures

We programmed the survey in Qualtrics and recruited 1000 Kenyan respondents from a consumer panel of Geopoll, implementing quotas for gender, age, and county of residency.<sup>5</sup> On average, respondents took 22 min (median = 15) to complete the entire survey, and each participant received a completion payment of 250 KES (2.20 USD at the time of the experiment), in addition to any eventual incentive payments.

## 5. Results

We randomly allocated 1000 participants to the four treatments, which resulted in 256 individuals in Control, 259 in Tips, 246 in Incentives, and 239 in Tips and Incentives. Individual characteristics are balanced across treatments (see Table A2).<sup>6</sup>

### 5.1. Descriptive statistics

Due to our quotas, half of our sample is female, 32% between 18 and 24 years, 27% between 25 and 34 years, and 41% 35 years and above. This implies that with 32 years on average our sample is older than the general Kenyan population but relatively comparable to the adult population (see Table A3). While respondents come from all over Kenya and are representative in terms of residency at the county level, urban participants are over-represented (50% as compared to 31% of the population in urban areas). Table A4 presents further descriptive statistics: Our sample is comparatively well-educated (73% have a post-secondary education), 78% self-classify as low-income and 36% have formal employment. As the design of the survey requires access to internet, it is not surprising that 99% have internet access and use social media on their phone. Almost all participants (96%) have recently used DFS on their phone and on average, participants use five different services with the most frequent ones being sending and receiving mobile money (89%), paying bills (71%), and conducting transactions involving an agent (55%).

In our sample, 96% report that they have been contacted by a scammer in the past.<sup>7</sup> Of those, 14% state having been contacted in the past week. The most common way of contact is reported to be SMS, followed by phone calls. Consistently with our findings from the social media and qualitative analysis (see Appendix C), the top three asks by scammers were to send money, to reverse a payment, and to share personal information. More than half of our sample report having ever been victimized.

### 5.2. Scam identification ability

We first present descriptive statistics from block 1, i.e., prior to the tips treatment. On average, participants correctly identified 71% of the six messages. Panel A in Figure A2 illustrates the distribution of SIA. Only 12% of all respondents correctly identified all six messages. Participants can make two kinds of identification mistakes: They might misclassify a scam (as a non-scam message), or they might misclassify a non-scam message (as a scam). On average, individuals classified 74% of scams and 66% of non-scams correctly. Confidence in SIA is high on average, at 4.23 out of 5 in block 1. Seventeen percent of participants always indicate the highest confidence score (see Panel B in Figure A2). SIA and confidence are positively correlated (Spearman’s rho=0.179,  $p < 0.001$ ).

Table 1 shows the correlates of SIA and confidence in block 1. Gender is the most robust and significant correlate of both SIA and confidence, with women having a 3 percentage point lower SIA score (equivalent to classifying 0.2 fewer messages correctly) and being less

<sup>6</sup> The data collection proceeded as planned and there were no changes to the pre-registered experimental design. In a few instances, we deviate from the pre-analysis plan, mostly for expositional clarity. We discuss all these changes in Appendix E.

<sup>7</sup> These numbers are substantially higher than the ones reported in the phone survey by Blackmon et al. (2021). This may be explained by several differences. First, in our survey, we provide participants with visual examples of scams that might make recall easier. Second, our sample is more educated than theirs, and they find that reports of scam contacts are positively correlated with education. Third, if reporting is influenced by social image concerns, online survey mode might increase reporting rates.

<sup>5</sup> For more detail on the recruitment strategy, see Appendix section C.3.

**Table 1**  
Correlates of scam identification ability and confidence.

	SIA			Confidence in SIA		
	(1)	(2)	(3)	(1)	(2)	(3)
<b>Demographics:</b>						
Female	-0.03*** (0.01)	-0.03*** (0.01)	-0.03** (0.01)	-0.11*** (0.04)	-0.10*** (0.04)	-0.10** (0.04)
Age in years	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.01** (0.00)	0.01** (0.00)	0.01** (0.00)
Post-Secondary education	0.03* (0.01)	0.02 (0.01)	0.02 (0.01)	0.11** (0.05)	0.10* (0.05)	0.09* (0.05)
Low income	0.02 (0.01)	0.02 (0.01)	0.02 (0.01)	0.04 (0.05)	0.03 (0.05)	0.03 (0.05)
Formal employment	-0.00 (0.01)	-0.01 (0.01)	-0.00 (0.01)	0.06 (0.04)	0.04 (0.04)	0.05 (0.04)
<b>DFS Use:</b>						
Low trust in DFS		0.01 (0.01)	0.01 (0.01)		-0.11** (0.05)	-0.11** (0.05)
High use of different DFS		0.03** (0.01)	0.03** (0.01)		0.05 (0.04)	0.05 (0.04)
<b>Scam Experience:</b>						
Contacted less than 1 week ago			-0.01 (0.02)			0.00 (0.06)
Victim of a scammer			-0.00 (0.01)			-0.02 (0.04)
N	997	997	991	997	997	991
R-Squared	0.05	0.05	0.05	0.03	0.04	0.03

Notes: Dependent variables are the share of correctly identified messages (SIA) in block 1 and average confidence ratings in block 1. *Female*, *Post-Secondary Education*, *Formal Employment*, *Low Trust in DFS*, *Contacted less than 1 week ago*, and *victim of a scammer* are binary indicators, *Low Income* and *High use of different DFS* are binary indicators for median splits. All variables rely on self-reports. All specifications control for the order of the two blocks and failing the attention check. The displayed coefficients are from OLS regressions. Robust standard errors are in parenthesis.

\*\*\*Indicate that the estimate is statistically significant at the 1% level.

\*\*Indicate that the estimate is statistically significant at the 5% level.

\*Indicate that the estimate is statistically significant at the 10% level.

confident in their ability. These results are consistent with the well-documented gender gap in financial literacy (Lusardi and Mitchell, 2014). Other demographic characteristics are at most weakly correlated with SIA. Age and having more than secondary education are positively correlated with confidence. Those who use a larger variety of DFS have a 3 percentage point better SIA score (they classify 0.2 more messages correctly). Low trust in DFS is associated with lower confidence. We find no significant association between individuals' scam experience (i.e., being contacted or victimized) and SIA or confidence.

Lastly, we assess the effect of incentives in block 1 on our four main outcome variables: SIA (the share of correctly identified messages), the share of correctly identified scams, the share of correctly identified non-scams, and the confidence level. Panel 1 in A5 shows that incentives have no significant effect on any of the outcomes. While we control for the incentive treatment in all the following analyses, we will focus on the two tips treatments for ease of exposition.

### 5.3. Effects of scam education

To test the null hypotheses that (i) tips (unincentivized) and (ii) tips (incentivized) have no effect on our main outcome variables, we estimate the following model:

$$y_i = \alpha_0 + \alpha_1 Tips_i^U + \alpha_2 Tips_i^I + \gamma_1 y_{0i} + X_i' \gamma_2 + Other_i \delta + \epsilon_i,$$

where  $y_i$  is our outcome variable measured in block 2.  $Tips_i^U$  indicates that individual  $i$  received the tips treatment without the incentives.  $Tips_i^I$  indicates that individual  $i$  received both the tips and incentives.  $y_{0i}$  controls for the baseline levels of the outcome variable from the first block.  $X_i$  is a set of individual characteristics for respondent  $i$ . These include gender, age, income, and education level.  $Other_i$  captures additional controls, such as the order of the two blocks and whether individual  $i$  received incentives (with no tips). We use robust standard errors  $\epsilon_i$ . Our coefficients of interest are  $\alpha_1$  and  $\alpha_2$ , i.e., the effect of tips without the incentives and with the incentives, respectively.

Column 1 of Panel 1 in Table 2 shows that tips do not increase SIA relative to the control group (no tips and no incentives). The same holds for tips with incentives. Columns 2 and 3 help explain why tips have no overall effect. While tips are helpful in increasing the share of correctly identified scams (Column 2), they decrease the share of correctly identified non-scams (Column 3). These effects do not depend on incentives.

Columns 4 to 6 present the treatment effects on confidence. Column 4 shows that, on average, individuals who received tips become more confident in their classifications. This increase is driven by participants becoming more confident in the classification of scams (Column 5). In contrast, the confidence in the classification of non-scams does not change with tips (Column 6), despite the worse performance (Column 3).<sup>8</sup>

Panel 2 in Table 2 shows the effect of our treatments on secondary outcomes. First, we find no significant effect of tips on trust in digital financial services.<sup>9</sup> Tips increase the time participants spend on the classification task in comparison to the control group only when the incentives are provided. Note, however, that we cannot statistically distinguish the effect of tips with and tips without the incentives. The former may induce higher effort (proxied by longer response times), but this does not lead to better outcomes. The last two columns show treatment effects on classifying all scams and all non-scam messages correctly, confirming the results from Panel A. Our results are not

<sup>8</sup> These averages might mask substantial heterogeneity. We hence try to assess to what extent changes in confidence coincide with improvements in SIA. Table A6 provides suggestive evidence that, on average, increases in confidence occur together with increases in SIA (Column 1). This association of SIA and confidence is particularly strong for scams (Column 2), but reversed for non-scams (Column 3): confidence does not increase while performance decreases.

<sup>9</sup> Note that we measure trust in DFS only once after all messages have been classified. We hence cannot control for a baseline level of trust.

**Table 2**  
Treatment effects.

Panel 1: Main outcomes

	Correctly identified messages			Confidence		
	SIA	Scams	Non-scams	SIA	Scams	Non-scams
Tips (unincentivized)	0.02 (0.02)	0.08*** (0.02)	-0.09*** (0.03)	0.12*** (0.04)	0.16*** (0.05)	0.06 (0.05)
Tips (incentivized)	0.03* (0.02)	0.08*** (0.02)	-0.07** (0.03)	0.08* (0.04)	0.08 (0.05)	0.07 (0.06)
Control Mean	0.70	0.69	0.71	4.20	4.13	4.33
p-value ( $Tips^U = Tips^I$ )	0.69	0.82	0.60	0.37	0.14	0.85
N	991	991	991	991	991	991
R-Squared	0.04	0.11	0.16	0.47	0.40	0.27

Panel 2: Secondary outcomes

	Low Trust in DFS	Response time SIA	All scams identified	All non-scam identified
Tips (unincentivized)	-0.01 (0.04)	0.11 (0.08)	0.10** (0.04)	-0.11** (0.04)
Tips (incentivized)	-0.02 (0.04)	0.21** (0.09)	0.11*** (0.04)	-0.08* (0.04)
Control Mean	0.32	2.21	0.30	0.52
p-value ( $Tips^U = Tips^I$ )	0.92	0.27	0.82	0.48
N	991	991	991	991
R-Squared	0.04	0.34	0.07	0.11

Notes: In Panel 1, the dependent variables are the share of correctly identified messages (SIA) in block 2, the share of correctly identified scams in block 2, the share of correctly identified non-scams in block 2, and the average confidence ratings in block 2 for all messages (confidence in SIA), for the scam messages, and for the non-scam messages. In Panel 2, the dependent variables are a binary indicator for low trust in DFS, the time spent on SIA in block 2, a binary indicator for classifying all scams correctly in block 2, and a binary indicator for classifying all non-scams correctly in block 2. All specifications include an indicator for the incentives treatment, the value of the outcome variable in block 1 (except for trust, which was only measured after block 2), and the full set of controls, i.e., variables displayed in Table 1 (female, age, post-secondary education, low income, formal employment, low trust in DFS (except for the effect on trust), above median use of different DFS, contacted less than one week ago, victim of a scammer), as well as indicators for the order of the two blocks and failing the attention check.  $Tips^U$  and  $Tips^I$  refer to Tips (unincentivized) and Tips (incentivized), respectively. The displayed coefficients are from OLS regressions. Robust standard errors are in parenthesis.

\*\*\*Indicate that the estimate is statistically significant at the 1% level.

\*\*Indicate that the estimate is statistically significant at the 5% level.

\*Indicate that the estimate is statistically significant at the 10% level.

driven by a lack of attention or a specific set of control variables (see Appendix B).

#### 5.4. Heterogeneity

We investigate who benefits from tips. Specifically, we explore treatment effects for respondents separately by the following characteristics: gender, age, education, income level, rural and urban areas as well as experience with DFS and scams. Figs. 2(a) and 2(b) plot the coefficients of SIA and confidence, respectively, for each subgroup.

First, we note that the directions of effects in most subgroups are consistent with our main results and most subgroups react equally to the tips treatments. In terms of SIA, tips, irrespective of the presence of incentives, appear to work better for those with post-secondary education and a more diverse use of DFS (using 5 or more different services). Recall that those with more DFS experience are also better at identifying scams in the baseline (see Table 1). This suggests that tips further increase the gap in SIA between inexperienced and experienced DFS users. Confidence increases for most subgroups, generally with only subtle differences between groups.

## 6. Discussion

We find no significant average effect of tips on SIA, but differential effects of tips on scams and non-scams. In this section, we present potential explanations for these results and underlying mechanisms. Additionally, we discuss how to interpret our effect sizes. Note that this section is exploratory in nature.

### 6.1. Exploring the effect of tips on SIA

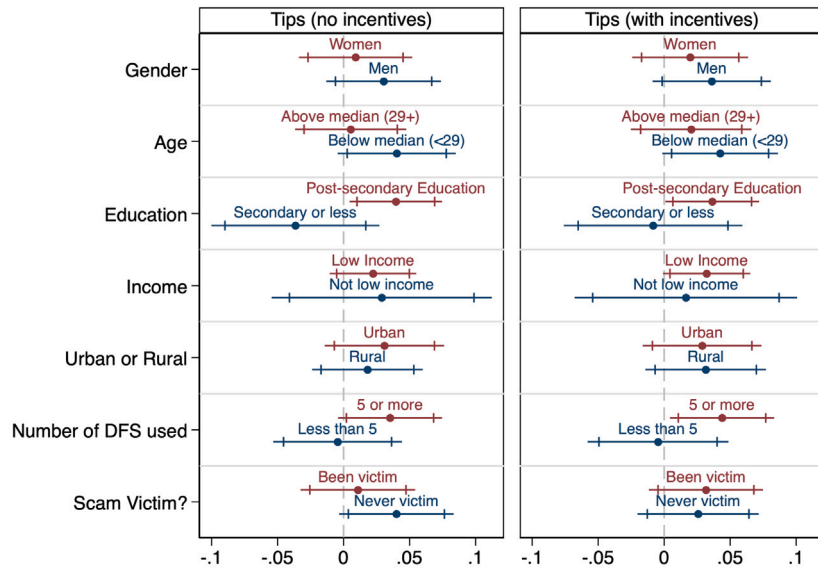
Our light-touch scam education in the form of scam tips does not improve SIA. However, tips improve the identification of scams, while they worsen the identification of non-scams. This pattern could emerge due to two reasons. First, tips may increase caution, such that participants are more likely to classify any given message as a scam. Second, not only scams but also non-scam messages may contain “scam markers”, such that tips “apply” to both scam and non-scam messages. In the former case, policymakers may want to weigh the benefits of improved scam identification against the costs of heightened classification mistakes for genuine communication — providing tips could still be welfare-improving if the cost of avoiding genuine communication is relatively low. In the latter case, it should be discussed whether tips can be refined and whether official communication can distinguish itself better from scams.

We analyze the effects of our treatments at the vignette level to shed light on potential mechanisms. To account for the fact that not all tips are helpful for all vignettes, we construct an indicator,  $ScamMarker_m$ , which captures whether at least one of the tips is helpful for correctly identifying the message as a scam. Only one scam message does not contain a scam marker while the other seven do. Yet, two out of the four official messages also contain a scam marker making them look like scams.

Fig. 3 plots the average marginal effects obtained from our estimates for the control group and the tips without incentives treatment in block 2.<sup>10</sup> In the left panel, we include all messages, in the center

<sup>10</sup> We focus on this comparison for ease of exposition; the effects for the tips with incentives treatment are qualitatively similar. We provide more detail in Appendix D.

(a) Scam Identification Ability (SIA)



(b) Confidence

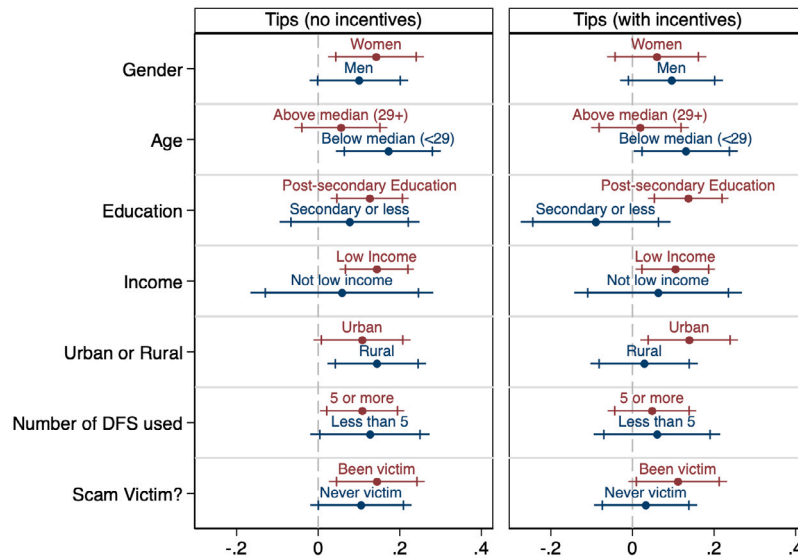


Fig. 2. Treatment effect heterogeneity.

Notes: Figures plot the OLS coefficients and the 90% and 95% confidence intervals from the estimating regressions in Panel 1, Table 2 (Column 1 for SIA and Column 4 for Confidence) separately for the different subcategories.

panel, we only include scam messages, and in the right panel non-scams.<sup>11</sup> Similar to our main results, we find no differential effect of our treatment on the share of correctly identified messages in block 2, irrespective of scam markers (left panel).

Focusing only on scams (center panel), tips significantly increase the share of correctly identified messages, independent of whether the message contains a scam marker or not. This is in line with the

<sup>11</sup> Note that the magnitudes cannot be compared across the panels as the share of correctly identified messages relies on six (left panel), four (center), and two messages (right panel), such that one mistake has a different magnitude in the three panels.

interpretation that participants become more cautious and hence more likely to classify any given message as scam when they receive tips. There is one caveat worth mentioning here. We only have one scam message without a scam marker. For non-scams, we see that tips do not increase the share of correctly identified messages for messages without a scam marker. However, if a scam marker is present, tips significantly *reduce* the share of correctly identified messages. This highlights the challenge of designing educational campaigns in a setting in which genuine communication contains scam markers.<sup>12</sup> We conclude that

<sup>12</sup> Note that scam markers in official communication are not specific to our experiment. Anecdotally, we were surprised to find other scam markers such

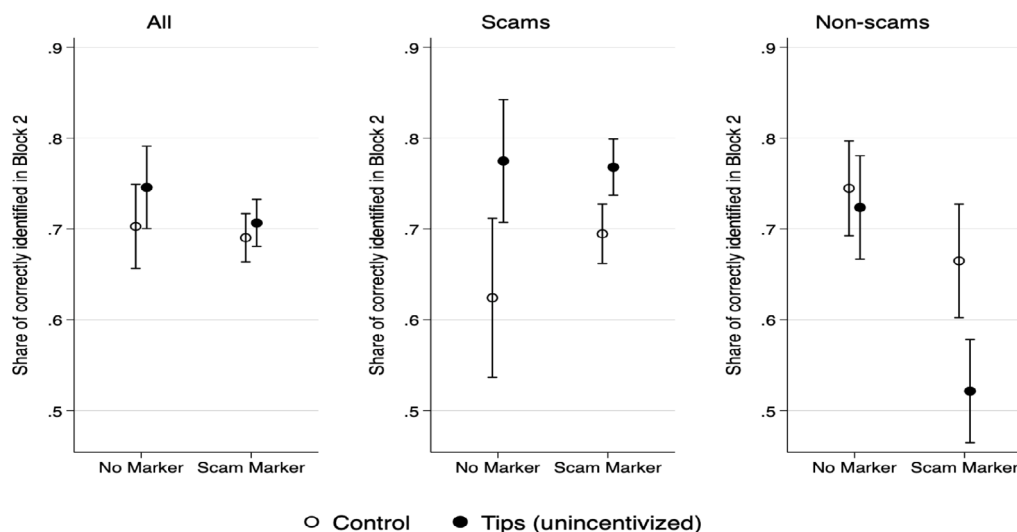


Fig. 3. Vignette-level effects by whether the message contains a scam marker.

Notes: Figures plot the average marginal effects of triple-differences estimation with 95% confidence intervals based on standard errors clustered at the respondent level (see also Appendix D). Scam Marker is an indicator for whether the message contains at least one of the scam markers the tips warn about. The left panel contains all vignettes, the center panel focuses on scams, and the right panel on non-scams. For ease of exposition, only the control and the Tips (unincentivized) treatment are displayed. The empirical specification contains the full set of interactions and demographic controls.

if non-scams can avoid scam markers, tips can be unambiguously beneficial in increasing scam detection irrespective of scam markers while not decreasing the correct classification of non-scams.

## 6.2. Interpretation of effect sizes

Our setting differs in several ways from the “real life”. For one, we abstract away from situational circumstances that may help classify messages. We also focus on messages that may be harder to classify than the average SMS individuals receive in Kenya. In general, without knowing all messages and the frequency at which they are being received, it is hard to interpret the absolute levels of our SIA measure. Thus, we mainly focus on differences in SIA between different groups, either defined by our treatments, or by demographics.

As to our treatment effects, we are primarily interested in their directions, and less so in the magnitudes. There are several reasons to believe that we estimate an upper bound of the effect of tips. First, our sample is literate and relatively educated and hence able to understand and apply the tips. In line with this, more educated and experienced DFS users appear to benefit more from the tips. Second, we provide tips when they are needed, in a more salient way than in real life. Additionally, as participants are aware they might face scams, they may pay more attention to tips than they would otherwise.

However, other points speak toward a lower bound of the effect. Being alert also means that the awareness-raising potential of tips is weakened, if not muted. As we find tips to improve the identification of scams even when attention is incentivized, this argument seems to have less bite. In addition, since we use common tips, participants may know them already. This is especially likely given that our sample is more educated and uses the internet more than the average Kenyan population. Finally, if average scams are less challenging to identify than our vignettes, we might estimate a lower bound, as the following analysis suggests. Using vignette-level data from block 1, we create a measure of difficulty and analyze treatment heterogeneity at the vignette level in block 2, analogously to the analysis of “scam markers”. For easy vignettes, we find a slight increase in SIA with tips for all

as urgency, all caps or shortened links in several of the official communication messages sent by banks and Safaricom.

messages, a positive and significant effect for scams, and no effect for non-scams. For difficult vignettes, tips significantly increase the correct classification of scams but significantly decrease the correct classification of non-scams (see Appendix D and Figure D2 for more details). Assuming that most official entities manage to communicate in easy-to-classify messages, we rather estimate a lower bound of the effectiveness of tips.

Lastly, we note that a limitation of our approach is that the effects of tips are examined using vignettes. While we make classification mistakes costly for half of our sample, this does not take into consideration that the costs of misclassifying scams and non-scams are likely different in practice. Real-life stakes could also be much higher than the experimental ones. Moreover, similar to other studies in the literature, we are not able to assess how SIA translates into fraud detection in practical settings and the likelihood of victimization (Burke et al., 2022). Our results suggest that tips can decrease classification errors for scams, but further testing and quantifying effects, also in terms of potential downsides for non-scam communication, remains an important question for further research.

## 7. Conclusion

We study a progressive DFS market in which phone scams are highly prevalent, develop a measure of scam identification ability, and experimentally test the effect of scam education in the form of tips. On average, we find no significant effect of tips on SIA. We explain this null effect by an increase in correctly identified scams, and a decrease in correctly identified genuine messages. Further analyses reveal that these differential effects appear to be driven by scam markers that are also present in some of the non-scam communication by banks or telecommunication companies. If such communication could be distinguished more easily from scams, tips on how to spot scams may have an unambiguously positive effect on SIA. Moreover, we show that tips lead to an increase in confidence, driven by higher confidence in classifying messages that are indeed scam. We also find suggestive evidence that tips do not make individuals overly confident. This is in line with specific subgroups, namely the more educated and more experienced, benefiting from the treatment and becoming more confident.

Our analyses reveal several reasons why scam tips, despite being a commonly used approach, might not be the silver bullet in addressing



the human factor in scam victimization. First, it is challenging, if not impossible, to provide tips that benefit all. Our findings suggest that tips, for example, benefit only the highly educated which potentially leads to a further increase in gaps between groups. Therefore, a more targeted approach may be necessary to reach everyone, and in particular, populations who may be more susceptible to scam victimization. Importantly, targeting is not only about the content, but also the medium used to educate consumers. For example, [Burke et al. \(2022\)](#) find that text-based messaging may work better for more educated populations, potentially explaining why our written texts work better for this subgroup. Second, it is difficult to communicate tips that apply to all kinds of scams. Tips in our setting seem to increase scam detection irrespective of scam markers, potentially due to an increase in caution. Moreover, as scams evolve dynamically, tips and guidance provided by authorities need to be revised regularly. Notifying consumers of these updates poses an additional challenge. Therefore, identifying new strategies for fraud prevention and scam awareness remains an important endeavor for future research.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Data availability

The data and replication files can be accessed at Harvard Dataverse: <https://doi.org/10.7910/DVN/HUYQZ>.

### Acknowledgments

For their valuable comments and suggestions, we thank seminar and conference participants at the ASFEE 2023, the CEPR WEFIDEV Workshop 2023, the DIW Finance and Development Workshop 2023, the GDEC 2023, IMEBESS 2023, the IPA Researcher Gathering 2021, the M-BEPS 2023, the University of East Anglia, the University of Essex, the University of Exeter, the University of Groningen, the Ludwig-Maximilian-University of Munich, and the University of Strathclyde. We are grateful to Lyne Chahed for her outstanding research assistance and to Nendo for excellent research support.

Funding from Innovations for Poverty Action Consumer Protection Research Initiative (BMG-19-10001-X5) is gratefully acknowledged. Jana Čahlíková acknowledges support from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy – EXC 2126/1-390838866. Eva Raiber acknowledges funding from the French government under the “France 2030” investment plan managed by the French National Research Agency (reference: ANR-17-EURE-0020) and from the Excellence Initiative of Aix-Marseille University - A\*MIDEX. Eva Raiber is a Research Affiliate at the Centre for Economic Policy Research (CEPR). Elif Kubilay is a Research Fellow at the IZA. Lisa Spantig is a CESifo affiliate. The study has received IRB approval from IPA (16014), and ethical approval from the University of Essex (ETH2021-1858), and is registered with the AEA registry (AEARCTR-0008754). All authors approved the final version of the manuscript. The order of the authors was randomized.

### Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jdeveco.2023.103147>.

### References

- Andersen, Jørgen Juel, Johannesen, Niels, Rijkers, Bob, 2022. Elite capture of foreign aid: evidence from offshore bank accounts. *J. Polit. Econ.* <http://dx.doi.org/10.1086/717455>.
- Annan, Francis, 2022a. Gender and financial misconduct: a field experiment on mobile money. <http://dx.doi.org/10.2139/ssrn.3534762>.
- Annan, Francis, 2022b. Misconduct and reputation under imperfect information. <http://dx.doi.org/10.2139/ssrn.3691376>.
- Balyuk, Tetyana, 2022. FinTech lending and bank credit access for consumers. *Manage. Sci.* <http://dx.doi.org/10.1287/mnsc.2022.4319>.
- Blackmon, William, Mazer, Rafe, Warren, Shana, 2021. *Kenya Consumer Protection in Digital Finance Survey Report*. Technical report.
- Brenner, Lukas, Meyll, Tobias, Stolper, Oscar, Walter, Andreas, 2020. Consumer fraud victimization and financial well-being. *J. Econ. Psychol.* 76, 102243. <http://dx.doi.org/10.1016/j.joep.2019.102243>.
- Burke, Jeremy, Kieffer, Christine, Mottola, Gary, Perez-Arce, Francisco, 2022. Can educational interventions reduce susceptibility to financial fraud? *J. Econ. Behav. Organ.* 198, 250–266. <http://dx.doi.org/10.1016/j.jebo.2022.03.028>.
- Chen, Yan, YeckehZaare, Iman, Zhang, Ark Fangzhou, 2018. Real or bogus: predicting susceptibility to phishing with economic experiments. *PLOS ONE* 13 (6), e0198213. <http://dx.doi.org/10.1371/journal.pone.0198213>.
- DeLiema, Marguerite, Deevy, Martha, Lusardi, Annamaria, Mitchell, Olivia S., 2020. Financial fraud among older Americans: evidence and implications. *J. Gerontol. Ser. B Psychol. Sci. Soc. Sci.* 75 (4), 861–868. <http://dx.doi.org/10.1093/geronb/gby151>.
- Dustan, Andrew, Hernandez-Agramonte, Juan Manuel, Maldonado, Stanislao, 2023. Motivating bureaucrats with behavioral insights when state capacity is weak: Evidence from large-scale field experiments in Peru. *J. Dev. Econ.* 160, 102995. <http://dx.doi.org/10.1016/j.jdeveco.2022.102995>.
- Engels, Christian, Kumar, Kamlesh, Philip, Dennis, 2020. Financial literacy and fraud detection. *Eur. J. Finance* 26 (4–5), 420–442. <http://dx.doi.org/10.1080/1351847X.2019.1646666>.
- Ensminger, Jean, Leder-Luis, Jetson, 2022. *Detecting Fraud in Development Aid*. Working Paper Series, 30768, National Bureau of Economic Research, <http://dx.doi.org/10.3386/w30768>.
- Fabregas, Raissa, Kremer, Michael, Schilbach, Frank, 2019. Realizing the potential of digital development: The case of agricultural advice. *Science* 366 (6471), eaay3038. <http://dx.doi.org/10.1126/science.aay3038>.
- Financial Institution Regulatory Authority, Investor Education Foundation, 2015. *The Non-Traditional Costs of Financial Fraud: Report of Survey Findings*. Technical report, Applied Research and Consulting, New York, NY.
- Fu, Jonathan, Mishra, Mrinal, 2022. Fintech in the time of COVID-19: Technological adoption during crises. *J. Financial Intermed.* 50, 100945. <http://dx.doi.org/10.1016/j.jfi.2021.100945>.
- Garz, Seth, Giné, Xavier, Karlan, Dean, Mazer, Rafe, Sanford, Caitlin, Zinman, Jonathan, 2021. Consumer protection for financial inclusion in low- and middle-income countries: Bridging regulator and academic perspectives. *Annu. Rev. Finan. Econ.* 13 (1), 219–246. <http://dx.doi.org/10.1146/annurev-financial-071020-012008>.
- Guiso, Luigi, Sapienza, Paola, Zingales, Luigi, 2008. Trusting the stock market. *J. Finance* 63 (6), 2557–2600. <http://dx.doi.org/10.1111/j.1540-6261.2008.01408.x>.
- Gurun, Umit G., Stoffman, Noah, Yonker, Scott E., 2017. Trust busting: The effect of fraud on investor behavior. *Rev. Financ. Stud.* 31 (4), 1341–1376. <http://dx.doi.org/10.1093/rfs/hhx058>, arXiv:<https://academic.oup.com/rfs/article-pdf/31/4/1341/29005106/hhx058.pdf>.
- He, Daixin, Lu, Fangwen, Yang, Jianan, 2023. Impact of self- or social-regarding health messages: Experimental evidence based on antibiotics purchases. *J. Dev. Econ.* 163, 103056. <http://dx.doi.org/10.1016/j.jdeveco.2023.103056>.
- Holst, Christine, Isabwe, Ghislain Maurice Norbert, Sukums, Felix, Ngowi, Helena, Kajuna, Flora, Radovanović, Danica, Mansour, Wisam, Mwakapeje, Elibariki, Cardellichio, Peter, Ngowi, Bernard, Noll, Josef, Winkler, Andrea Sylvia, 2021. Development of digital health messages for rural populations in Tanzania: Multi-and interdisciplinary approach. *JMIR MHealth UHealth* 9 (9), e25558. <http://dx.doi.org/10.2196/25558>.
- Johnson, Eric J., Meier, Stephan, Toubia, Olivier, 2019. What's the catch? suspicion of bank motives and sluggish refinancing. *Rev. Financ. Stud.* 32 (2), 467–495. <http://dx.doi.org/10.1093/rfs/hhy061>.
- Koyama, Naoko, Totapally, Swetha, Goyal, Shruti, Sonderegger, Petra, Rao, Priti, Gosselt, Jasper, 2021. *Kenya's Digital Economy: A People's Perspective*. Technical report.
- Lusardi, Annamaria, Mitchell, Olivia S., 2014. The economic importance of financial literacy: Theory and evidence. *J. Econ. Lit.* 52 (1), 5–44. <http://dx.doi.org/10.1257/jel.52.1.5>.
- McAlvanah, Patrick, Anderson, Keith, Letzler, Robert, Mountjoy, Jack, 2015. *Fraudulent Advertising Susceptibility: An Experimental Approach*. Technical report.
- Moustafa, Ahmed A., Bello, Abubakar, Maurushat, Alana, 2021. The role of user behaviour in improving cyber security management. *Front. Psychol.* 12.
- Norris, Gareth, Brookes, Alexandra, Dowell, David, 2019. The psychology of internet fraud victimisation: A systematic review. *J. Police Crim. Psychol.* 34 (3), 231–245. <http://dx.doi.org/10.1007/s11896-019-09334-5>.

Pazarbasioglu, Ceyla, Mora, Alfonso Garcia, Uttamchandani, Mahesh, Natarajan, Harish, Feyen, Erik, Saal, Mathew, 2020. Digital financial services. In: World Bank Symposium. p. 54.

Scheibe, Susanne, Notthoff, Nanna, Menkin, Josephine, Ross, Lee, Shadel, Doug, Deevy, Martha, Carstensen, Laura L., 2014. Forewarning reduces fraud susceptibility in vulnerable consumers. *Basic Appl. Soc. Psychol.* 36 (3), 272–279. <http://dx.doi.org/10.1080/01973533.2014.903844>.

Sheng, Steve, Magnien, Bryant, Kumaraguru, Ponnurangam, Acquisti, Alessandro, Cranor, Lorrie Faith, Hong, Jason, Nunge, Elizabeth, 2007. Anti-phishing phil: the design and evaluation of a game that teaches people not to fall for Phish. In: Proceedings of the 3rd Symposium on Usable Privacy and Security. In: SOUPS '07, Association for Computing Machinery, New York, NY, USA, pp. 88–99. <http://dx.doi.org/10.1145/1280680.1280692>.