# A Multilingual Comparison of Email Scams

Duo Pan
*Pennsylvania State University*

Ellen Poplavska
*Pennsylvania State University*

Yichen Yu
*Pennsylvania State University*

Susan Strauss
*Pennsylvania State University*

Shomir Wilson
*Pennsylvania State University*

## Abstract

Email scams threaten the security of all email users, regardless of their language. Scam detection methods using email origins and denylists remain imperfect, and user education about scam emails is more necessary than ever. However, existing research and educational content on email scams focus on scams written in English, leaving open the possibility that scams may differ in other languages. To determine whether cultural gaps exist between email scams in different languages, we analyzed scam emails in English, French, and Russian collected from Anti-Fraud International, a popular anti-scam web forum. In this paper, we will examine similarities and differences between these datasets, revealing a need for a more culturally-aware approach to user scam education.

## 1 Introduction

A scam email is an email sent for a fraudulent purposes. The primary purpose of these emails is to steal email users' money or identity [7]. Scams are a prevalent and costly security risk faced by email users. Phishing emails, one type of scam, cause the loss of $500 million each year in the United States [3].

Although email scams are a global problem, current scam email research is heavily English-centric [6]. User education methods play a crucial role in battling scam emails, and cultural differences should be considered to create materials that are useful for non-English speakers. Because of this gap in existing research, we have compiled a multilingual scam email dataset and analyzed it from computational and linguistic perspectives. We built our dataset of English, French, and Russian

scam emails from the Anti-Fraud International online forum, where users post scam emails they have received in different language sub-forums. Because these results come from a dataset of self-reported scam emails on a primarily English website, it is possible that the trends we present may not represent all scam emails. After constructing these datasets, we explored the most common subjects of scam emails through a system of categories and identified the most frequent words in scam emails in each language. We also analyzed greetings and sender IP addresses. We further examined how scammers use word choice to build trust with email respondents within each language corpus. Finally, we sought to identify how scammers manipulate culturally different audiences by analyzing frequent words and themes across languages.

## 2 Related Work

The current defense methods against scam emails can be divided into three categories. The first is a denylist mechanism, which extracts sender email addresses and adds them to a list. If an email address has been added to this denylist, email sent from it is blocked by filters. Denylists are updated based on user reports of scams. Since phishing is one category of scam emails, Sheng et al. [8] did an empirical analysis of phishing email denylists. Although the mechanism can detect many phishing emails, it always fails to detect new scams, because the new email address is not in the previous list. These drawbacks call for the automatic detection of scams.

The second method is machine learning, in which researchers consider scam email detection a classification problem. To detect scams using machine learning, researchers extract linguistic features from scams and legitimate emails, then apply machine learning classification models to detect scams. Harikrishnan et al. [5] found that the Decision Tree and Random Forest machine learning methods achieved the highest accuracy when training algorithms to detect scams.

The third method of defending against scam emails and preventing loss is user education. Researchers Almomani et al. [1] have tried to improve users' awareness of scam emails

with online training and testing. Governments and non-profit organizations post online information and training modules to help users avoid scam attempts. Diaz et al. [2] completed a study of user scam susceptibility in academia and found that, counterintuitively, students who had received phishing awareness training were more susceptible to phishing links. Only 28% of students without phishing training clicked the phishing links, while 42% of students with phishing detection training clicked on phishing links. This signals that the user education paradigm may be inadequate. Therefore, more culturally-aware user education approaches are necessary.

Researchers have also studied social engineering methods in scam emails. Markus Jakobsson [6] found that the use of targeted scam emails aimed at a specific group of people had increased from 2006 to 2014, while non-targeted scam emails decreased. Jakobsson also analyzed the various principles of persuasion that scammers used, such as authority, social proof, and commitment. Finally, Jakobsson completed case studies of sales, romance, and business scams by exploring the scam activities of each. Jakobsson used honeypot ads to collect scammer information and attributed 50% of the sales scam attempts on Craigslist to just ten groups of scammers.

Christopher J. Hadnagy and Michele Fincher [4] focused on the victims of scam emails, exploring how victims made the decision to open and interact with scams. Hadnagy and Fincher provided insights for user education methods by analyzing the scam email problem from the victim's perspective.

## 3 Corpus Creation

### 3.1 Sources for Scam Emails

We collected English, French, Russian scam emails from the Anti-Fraud International online forum[1]. We created a system of seven categories describing common types of scam emails, which can be applied to emails in all languages. Then, we selected the top three topics within English, French, and Russian forums and used a web crawler to collect scam emails. We will discuss more details and provide Table 1, which shows scam emails we crawled, in the Datasets section of this paper.

### 3.2 Annotation Scheme

The first annotation scheme is a list of seven mutually exclusive categories sorting the forums into different scam types across all languages. The second is a list of words frequently found in scam emails about different topics across languages. For French and Russian scam emails, we extracted frequent word lists from both the original and Google-translated corpora and analyzed the context of each word. This second scheme will be discussed further in the Analysis section. The use of machine translation is a limitation of our work; although language experts on our team spot-checked some

<hr>

[1]https://antifraudintl.org/

translations, a comprehensive audit of the translations was infeasible.

We built a corpus of scam emails in different languages from the Anti-Fraud International forums, where users post scam emails they have received to educate visitors of the website. Each of fifteen language forums on the site is divided further into a number of topic sub-forums. Rather than individually sorting all of the emails, we sorted the topic sub-forums on the language forums that contained complete usable scam email content into seven mutually exclusive categories: helping, profit, transaction, phishing, extortion, romance, and miscellaneous. In order to place each topic into a category, we read a sample of its emails and considered the title, guidelines, and description provided by the administrators of the forum.

The forum title keywords that guided forum sorting decisions are listed in Figure 1. All words are in English, either originally or Google-translated. Each category in this figure is split further into broad themes. These are not utilized for quantitative analysis in this paper, but are an opportunity for future study. The size of the boxes does not indicate the number of emails or threads in each category. Table 1 displays the number of forum threads in each category in each language.

The helping category includes all emails with a primary appeal to the recipient's generosity, asking them to help someone for little or no compensation. The profit category encompasses emails that appeal to the recipient by promising some benefit to the receiver. The transaction category includes scams centered around two main types of transactions: selling products or services to the recipient and offering work to the recipient. The phishing category includes scams with the primary intention of harvesting the recipient's personal information. The extortion category encompasses emails that primarily use fear to motivate recipients to provide the requested information or funds. The romance category includes emails written as if to begin a romantic relationship. The miscellaneous category, which made up a minority of all corpora, was necessary because each forum had at least one miscellaneous sub-forum.

### 3.3 Datasets

Based on the annotation scheme described above, we built a dataset of scam emails from the online forum. Figure 2 shows the pipeline from the original forum to the corpora we built.

First, we determined the number of forum threads in each of the seven categories within each language. Next, we selected the three categories with the most forum threads. Thread count is not an exact equivalent to precise email count, but most threads contain one scam email, so thread count provides an approximation for the sizes of the scam email corpora. Using thread counts in each category in each language, we found the top three categories in English, French, and Russian. Then, we selected the topic sub-forum in each category of each language with the greatest number of threads as a representative of the category. Table 2 summarizes the scam emails crawled
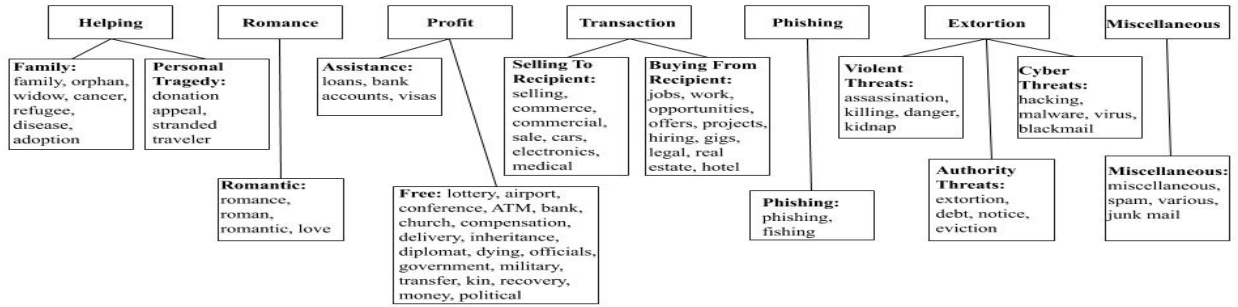
Figure 1: An illustration of all forum title keywords and how they fall into the category system

| Number of Forum Threads in Each Category by Language | | | | | | |
|---|---|---|---|---|---|---|
| Language | Helping | Romance | Profit | Transaction | Phishing | Extortion | Miscellaneous |
| English | 5433 | 3400 | 64456 | 19416 | 1100 | 442 | 201 |
| French | 217 | 575 | 1944 | 198 | 90 | 0 | 104 |
| Russian | 14 | 27 | 141 | 0 | 0 | 0 | 6 |

Table 1: Quantities of forum threads in each scam category by language, not including emails from personal inboxes
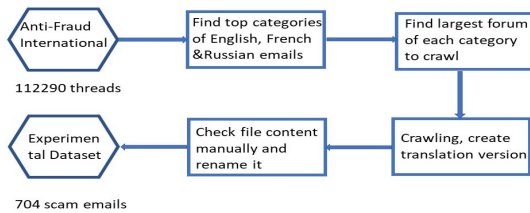


Figure 2: The pipeline from online forum to dataset

| Number of Crawled Scam Emails in Each Language | | | |
|---|---|---|---|
| Language | Top Categories | Top Forums | Number |
| English | Profit | Government | 98 |
| | Transaction | Business | 99 |
| | Helping | Orphans | 97 |
| French | Profit | Dying | 100 |
| | Romance | Romantic | 99 |
| | Helping | Orphans | 100 |
| Russian | Profit | Lottery | 70 |
| | Romance | Romantic | 27 |
| | Helping | Cancer Victims | 14 |

Table 2: Number of Crawled Scam Emails in Each Language

to represent English, French and Russian. It lists the three top categories with the most emails within each language, as well as the topic sub-forum representing each category. For instance, the top category of Russian emails was profit. Of profit emails, most came from the sub-forum titled "lottery", which contained 70 of the 141 Russian profit emails.

Russian had a smaller number of scams than English and French, so we under-sampled English and French scams to balance the corpus. We then removed files that were not complete emails. Then we separated sender information from the email body and named each file with the email index, language, source, category, and version (original or Google-translated).

## 4 Analysis

In this section, we will present case studies across different languages. First, we will compare orphan scam emails in the English and French forums to represent helping scam emails.

Then, we will compare romance scam emails in the French and Russian forums to represent romance scams.

## 4.1 English and French Orphan Scam Emails

We compared English and French orphan scam emails from two perspectives. First, we counted the word frequency of 97 English orphan scam emails and 100 French orphan scam emails. After tokenization, punctuation removal, lemmatization, and part-of-speech tagging, we calculated the 100 most frequent words in both the original and Google-translated emails. We found 72 overlapping words between the 100 most frequent words in English scams and 100 most frequent words in French scams. This reveals high lexical similarities between English and French orphan scam emails.

We divided these 72 frequently-occurring words into three types. The first word type is relevant to family and disaster, and is used to create an empathetic scenario. A typical scenario might be that the sender became an orphan because of an earthquake and received a large inheritance, which they plan to invest with the receiver's help. Frequent words include father, mother, and death. The second word type describes scammers' request, or call to action. After depicting their situation, the scammer proposes what they want the email receiver to do. Action words such as "help," "reply," "contact," and "please" are used. The last word type invokes money or finance. For example, both the English and French corpora included "bank," million," "fund," "business," and "account." The scammers used these words to promise a reward, or to attempt to steal receivers' money or access their accounts.

In addition to similar frequent words, English and French scam emails have context similarities. For example, scammers often used vague greetings. The English scammers used "my dear," "hello dear," or "my beloved," while French scammers used "bonjour". None of these indicate the receiver's name.

There are also differences between the English and French corpora. The English scammers preferred to use Gmail accounts (66%), while the French scammers preferred Yahoo Mail (61%). The English scammers created more varied disaster scenarios than the French scammers. In 100 French orphan scam emails, 60 emails mentioned the political situation in Côte d'Ivoire (Ivory Coast), while the English scammers invoked earthquakes, and parents' heart attacks more often.

## 4.2 French and Russian Romance Scam Emails

We also compared the 100 most frequent words for French and Russian romance scams, finding 52 overlapping words, particularly about building relationships. These included "want," "friend," and "relationship." Both Gmail and Yahoo accounts were commonly used by French and Russian scammers.

However, French romance scams included some particular frequent locations, including Canada, London, and Cotonou.

In Russian romance scams, the most frequent location words were Senegal and Sonatel (the principal telecommunications provider in Senegal). This may be because 48% of scammers' IP addresses are from Senegal. According to Jakobsson's findings [6], most Russian romantic scams are "traditional" romantic scam emails, because they originate in West Africa. Unlike Russian scammers, 35% of the French romance scammers used Alice Italy as their telecommunications provider.

## 5 Discussion and Future Work

Our work revealed differences in topics between our language-specific datasets. Profit and helping scam emails are the most common scam categories for English, French, and Russian. However, each language also has a unique top category. The English corpus has more transaction scam emails, while the French and Russian corpora have more romance scams.

Even across the same category, emails in different languages diverge in content. For instance, the English orphan emails used more words appealing to sympathy. In the romance category, Russian romance emails were more similar to advance-fee scams than French emails, in that about half of Russian scammers' IP addresses are from West Africa.

Based on what we found above, we see the possibility of improving email user education. For example, based on the category distribution of each language, we can focus on different scam categories when we develop user education in different languages. We can provide more specific guidance according to our analysis. When we educate Russian-speaking users on how to avoid scam emails, for example, we can offer more details about advance-fee scams. It could also be beneficial for all user education efforts to include more diverse examples based around our seven category system, as many training modules are currently purely phishing-focused.

In the future, we plan to continue comparing categories of emails within our corpus and crawl more emails in other languages, such as Chinese and Korean, and identify more cultural differences. We plan to develop guidelines for culturally conscious user education methods. For example, educators may ask users to choose all languages they often use on the Internet, then customize scam email samples to improve users' awareness of the scams they might encounter.

## References

[1] Ammar Almomani, BB Gupta, Samer Atawneh, Andrew Meulenberg, and Eman Almomani. A survey of phishing

email filtering techniques. *IEEE communications surveys & tutorials*, 15(4):2070–2090, 2013.

[2] Alejandra Diaz, Alan T Sherman, and Anupam Joshi. Phishing in an academic community: A study of user susceptibility and behavior. *Cryptologia*, 44(1):53–67, 2020.

[3] Yong Fang, Cheng Zhang, Cheng Huang, Liang Liu, and Yue Yang. Phishing email detection using improved rcnn model with multilevel vectors and attention mechanism. *IEEE Access*, 7:56329–56340, 2019. https://ieeexplore.ieee.org/iel7/6287639/6514899/08701426.pdf.

[4] Christopher Hadnagy and Michele Fincher. *Phishing dark waters*. Wiley Online Library, 2015.

[5] NB Harikrishnan, R Vinayakumar, and KP Soman. A machine learning approach towards phishing email detec-tion. In *Proceedings of the Anti-Phishing Pilot at ACM International Workshop on Security and Privacy Analytics (IWSPA AP)*, volume 2013, pages 455–468, 2018.

[6] Markus Jakobsson. *Understanding social engineering based scams*. Springer, 2016.

[7] Mehrbod Sharifi, Eugene Fink, and Jaime G Carbonell. Detection of internet scam using logistic regression. In *2011 IEEE International Conference on Systems, Man, and Cybernetics*, pages 2168–2172. IEEE, 2011. https://kilthub.cmu.edu/articles/Detection_of_Internet_Scam_Using_Logistic_Regression/6473309/files/11902892.pdf.

[8] Steve Sheng, Brad Wardman, Gary Warner, Lorrie Cra-nor, Jason Hong, and Chengshan Zhang. An empirical analysis of phishing blacklists. 2009.