

A Study on Evolution of Email Spam Over Fifteen Years

De Wang, Danesh Irani, and Calton Pu
College of Computing
Georgia Institute of Technology
Atlanta, Georgia 30332-0765
Email: {wang6, danesh, calton}@cc.gatech.edu

Abstract—Email spam is a persistent problem, especially today, with the increasing dedication and sophistication of spammers. Even popular social media sites such as Facebook, Twitter, and Google Plus are not exempt from email spam as they all interface with email systems. With an “arms-race” between spammers and spam filter developers, spam has been continually changing over the years. In this paper, we analyze email spam trends on a dataset collected by the Spam Archive, which contains 5.1 million spam emails spread over 15 years (1998-2013). We use statistical analysis techniques on different headers in email messages (e.g. content type and length) and embedded items in message body (e.g. URL links and HTML attachments). Also, we investigate topic drift by applying topic modeling on the content of email spam. Moreover, we extract sender-to-receiver IP routing networks from email spam and perform network analysis on it. Our results show the dynamic nature of email spam over one and a half decades and demonstrate that the email spam business is not dying but changing to be more capricious.

Keywords—email; spam; evolution

I. INTRODUCTION

Email is used everyday as a method to communicate, both for individuals and businesses, but also as an information management tool [1]. What started primarily as a person-to-person communication medium has spread widely to one-to-many (e.g. mailing-lists) and many-to-one (e.g. forwarded traffic) communication medium [2]. As social media has grown dramatically, email also enhances the functionality provided by them. For instance, users are sometimes given pseudo-email addresses which can be used to receive email on the social network as well as email can sometimes be used to interact with the social network using specially crafted email addresses.

Spam is unsolicited and unrelated content sent to users, which most commonly is associated with email, but also applies to several different domains including instant messaging, websites, and Internet Telephony [6], [7], [8], [9], [10]. Spam degrades a user’s experience as, by definition, it is an annoyance and gets in the way of users consuming non-spam content. In an extreme case, spam can be seen as a denial of information preventing user’s from finding non-spam content.

In August 1998, Cranor et al. [3] described the rapidly growing onslaught of unwanted email and since then the volume of spam has grown even more as the amount of all email sent has grown exponentially. Constituting an annoyance, email spam

10% of overall mail volume in 1998, which results in an enormous burden on the thousands of email service providers (ESPs) and millions of end users on the Internet [5].

In addition to being on the receiving side of spam, ESPs need to invest in developing filters to combat the spammers and likewise spammers evolve to avoid spam filters. The co-evolution nature of spammers and spam filters is an “arms-race”, which has resulted in numerous publications employing adversarial strategies to tackle the spam problem [12], [13], [14]. Pu et al. [15] and Fawcett [16] developed techniques for characterization and measurement of email spam trends and researchers have also examined other types of spam including phishing [17] and Web spam [18]. In addition, Guerra et al. [19] compared the effectiveness of old and recent filters over old and recent spam to obtain spam trends on email spam dataset.

In this paper, we investigate the trends of email spam in terms of content, topics, and sender-receiver network over 15 years by performing an evolutionary study on the Spam Archive dataset [23]. We aim to answer the question of whether the email spam business is dying (also, as identified by our title). More concretely, we make the following contributions:

- First, we perform a long-term evolutionary study on a large email spam dataset, which includes statistical analysis, topic modeling and network analysis.
- Second, we demonstrate the changes of email spam over time with respect to contents and spammer behaviors.
- Lastly, we prove that email spam business is not dying but is becoming sophisticated.

The remainder of the paper is organized as follows. We motivate the problem further in Section II. Section III introduces the Spam Archive dataset used in our study. Section IV presents the analysis performed on the dataset and findings derived from the results. Section V discusses the future of email spam business and the limitations of our study. We talk about related work in Section VI and conclude the paper in Section VII.

II. MOTIVATION

The paper is inspired by an article by Kaspersky labs [20] named “The dying business of email spam” [21], which stated that “Spam email is on the wane. And no one on God’s green

Earth is going to miss it”. The conclusions were based on their annual report [22] citing that the share of spam in email traffic decreased steadily throughout 2012 to hit a five year low.

We are excited by the decline in the volume of email spam but it also raises the question as to whether the email spam business is dying and will continue to decline. Besides the volume change, we also consider the quality of email spam and the impact, which may be constituting a new trend of email spam business. For instance, spammers may post email spam in a more complicated way using spoofed email addresses and changing email relay servers. Those kind of email spam may slip away under the inspection of spam filters. Thus, it motivated us to investigate the evolution of email spam using advanced techniques such as topic modeling and network analysis. We try to find out the real trend of email spam business through email content, meta information such as headers, and sender-to-receiver network over a long period of time.

III. DATA COLLECTION

In this section, we introduce the Spam Archive dataset and show the overview of the dataset used in our study.

Spam Archive dataset [23] is collected by Bruce Guenter since early 1998 using honey-pot addresses. The project is still ongoing with monthly releases of new email spam. Since it provides a continuous long-term email spam data source from a consistent source, it is an excellent dataset for our investigation into spam trends. The volume of email messages received over the 15 years is shown in Fig. 1, with the date on the x-axis and log-scale volume of email messages received per month on the y-axis. From the figure we see that email spam volume grows steadily over time. For the spike of email spam during 2006, Bruce Guenter has attributed this to one of the spam traps having a wild-card address which received increasingly large amounts of spam which was subsequently disabled after 2006, since most of the spam was duplicates of other spam received.

Besides showing the trend of overall volume of email spam, we also present the volume changes monthly for different years in Fig. 2, with the month of the year on the x-axis

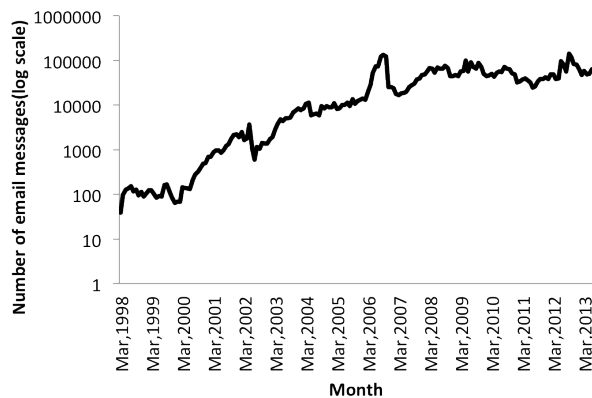


Fig. 1. Number of email messages (per month) over time

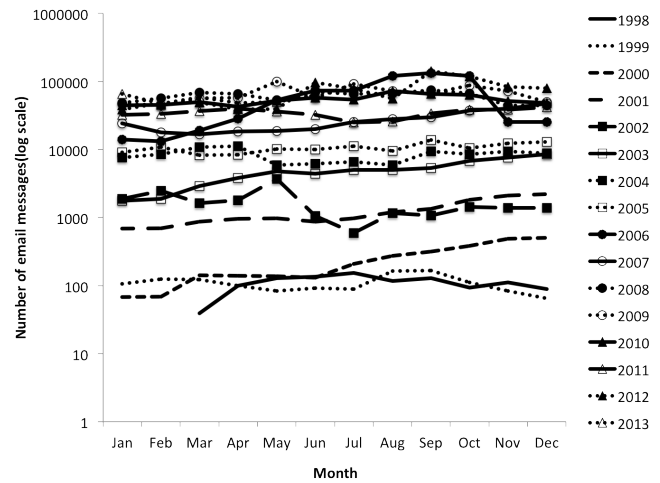


Fig. 2. Number of email messages in month order for different years

and the log-scale volume of spam messages per month on the y-axis. It shows volume trends over the previous 15 years. The volume of email spam is not always increasing over time such as the email spam volume changes during 1999. Some years’ volumes also shows fluctuations over time. For instance, during 2002, the volume first went up in May and decreased dramatically afterward until July. Several factors may have contributed to this change such as new strategies used by spammers (e.g. image spam is introduced in emails), improved spam filters (e.g. URL analysis tool is adopted) and even political influence from governments (e.g. Electronic Communications and Transactions Act, 2002 [24]). We investigate the details and potential reasons of these changes in more detail in the following sections.

IV. DATA ANALYSIS

In this section, we start with content analysis of Spam Archive dataset, followed by topic modeling and network analysis.

A. Content Analysis

The two main types of email message content are “Text” and “Multipart”. Messages in type “Text” are simple text messages while messages in type “Multipart” have parts arranged in a tree structure where the leaf nodes are any non-multipart content type and the non-leaf nodes are any of a variety of multipart types [25]. To have a better sense of the distribution of main types in email spam, we show the main type distribution in different years in Fig. 3.

Fig. 3 demonstrates that the distribution of two main types in our dataset changed over time. For instance, before 2003, more email spam had the message format in the main type “Text”. After that, the two main types almost occupied the same percentage until 2010. The new trend is that email spam is using more messages in main type “Text” (e.g. the percentage of email spam in main type “Text” is over 80% for the half year of 2013).

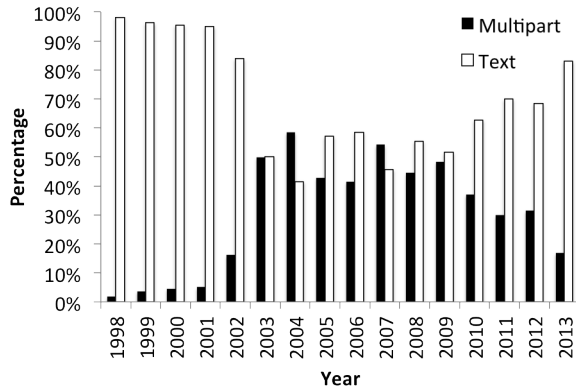


Fig. 3. The distribution of main types of message content

Next thing we are interested in is the embedded items in email spam such as HTML web page, images, and URL links. After scanning all email spam in our dataset, we present the distribution of embedded items in email spam over time in Fig. 4.

Fig. 4 shows that low percentage of email spam, which was always less than 5% in our dataset, contained image attachments. On the contrary, more email spam had embedded HTML web pages and URL links. But the percentages of email spam containing HTML web pages and URL links changed dramatically over time. Several peaks and valleys appeared over 15 years in the Fig. 4. For instance, HTML pages had peaks in 2003, 2007, and 2009 and valleys in 2006 and 2008. While for URL links, peaks appeared in 2004, 2008 and 2012 and valleys appeared in 2006 and 2011. Since HTML page normally carries URL links, they should have similar fluctuations along the time. However, we observe that an exception occurred after 2011. The percentage of email spam containing HTML web pages decreased suddenly after 2009. While the percentage of email spam containing URL links dropped down along with HTML web pages until 2011 and it increased sharply afterwards. One possible reason is that

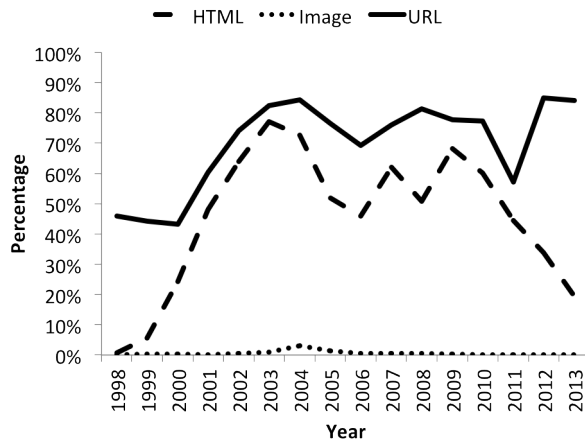


Fig. 4. The distribution of embedded items in email spam over time

more URL camouflage techniques, which are quite efficient in avoiding spam filters, appeared such as shortened URLs and hidden URLs in recent years. To investigate further the trend of URL links, we aggregate all URL links on a yearly basis for email spam that contain URL links and show the cumulative distribution of URL links in email spam in Fig. 5 (1998 – 2012). The data of 2013 is not included due to that it only contains half year data.

Fig. 5 shows the number of URL links for the majority of email spam is below 10. Only a small portion of email spam have more than 1,000 URL links which may be embedded in different depths of email messages. Even though the densities of URL links in email spam changed variously, email spam contained more and more URL links over time.

In addition to looking into embedded items, we also investigate the top n -grams in email spam over time. The tool we used for obtaining n -grams of email spam is Perl’s module Text::Ngrams [26]. First, we need to clean our dataset by filtering out stop words and stripping out HTML tags. And then we calculate top-10 n -grams (n ranges from 1 to 3) on a monthly basis over 15 years. Due to space limit, we only list the top-10 n -grams starting from June 1998 to June 2013, which is shown in Table I.

In Table I, $\langle N \rangle$ denotes any number sequence. Top-10 n -grams set contained different words or word sequences along the time, showing different topics as well. For instance, the n -grams set in June 1998 tells us that the email spam was advertising fake dental services using attractive words such as “free”, “nationwide near”, and “month save average”. The n -grams set in June 2003 was about marketing and market leaders leading people to click external URL links. The n -grams set in June 2008 was about DASS (Defensive Aids Sub System) [27] which is a fighter system from European countries. After checking the original email, it is a trap news or game to attract the email receivers to enter into. The n -grams set in June 2013 was more related to new media announcement and membership registration. Moreover, the differences indicate the topic drift in email spam over time (e.g. from fake advertising to fake registration services). To learn more about the topic drift of email spam, we will apply topic modeling on the dataset next.

B. Topic Modeling

Topic modeling is defined as a technique that looks for patterns in the use of words and it is an attempt to inject semantic meaning into vocabulary, in which a “topic” consists of a cluster of words that frequently occur together [28]. The tool we used in our topic modeling is a machine learning toolkit for language named “MALLET” [28]. It provides an efficient way to build up topic models based on Latent Dirichlet Allocation model (LDA) [29].

To simplify the illustration, we set up the number of topics to 10 in the data processing. After the calculation, we obtain the word (also called term) lists associated with topics and topic composition for different months over time, which is shown in Table III and Fig. 6.

In Table III, it shows the topic name and the samples of most related terms. After the topic modeling, we only have the word

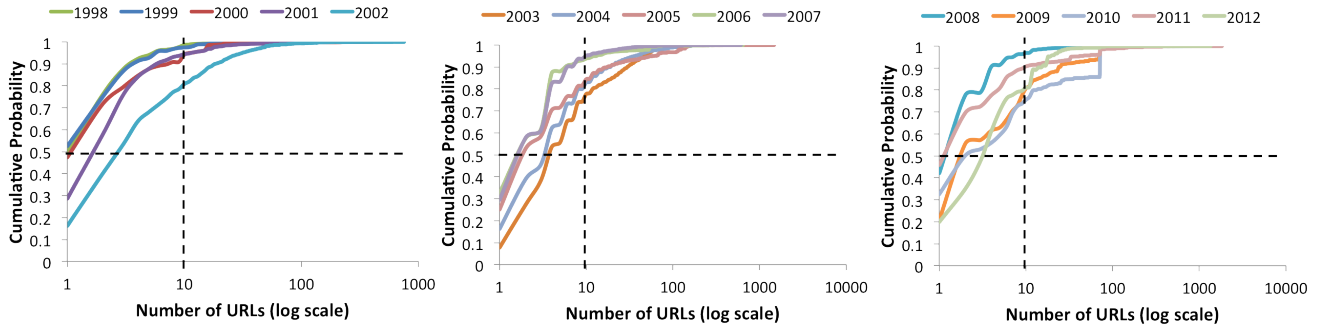


Fig. 5. Cumulative distribution of URL links in different years

TABLE I. List of top-10 n -grams every 5 years on a monthly basis (n ranges from 1 to 3)

June, 1998	June, 2003	June, 2008	June, 2013
dental free plan <N> details call please doctor dentistry procedures	<N> click email information bait mail free message work please	<N> euro dass online http mail original super time active	<N> important garden class email media screen dark right registration
plan free teeth whitening nationwide near waiting periods root canals details june dental procedures canals crowns doctor locator polishing fillings	<N> <N> email bait august <N> market information world leader auction records remove email reply message link work leader market	<N> <N> euro euro super active active euro tabs doses kinder dass autopilot dass original stress stress angst angst dass	<N> <N> garden <N> <N> garden <N> garden important media important important dark skin screen class class important rights reserved
sealants prevent cavities doctor locator number crowns dentures braces problems qualify waiting month save average receive optical plan call <N> please canals crowns dentures optical plan free plan receive optical	<N> <N> <N> world leader market leader market information case link work demander plus figurer allow mail removed removed thank operation modifier sera effective message modifier sera effective coop demander	<N> <N> <N> euro euro euro super active euro dass kinder dass dass autopilot dass dass dass kinder autopilot dass dass original stress angst kinder dass super	garden <N> garden <N> garden <N> <N> <N> <N> important media screen media screen class important important media class important media screen class important limited become member become member soon

or term clusters for each topic which has not been labeled. Based on associated terms with each topic and experience with email spam, we label the topics as “Account Information”, “Order Information”, “Business News”, “Sales News”, “Adult Product”, “Software Product”, “Official News”, “Free Product”, “Medical Product”, and “Newsletter” separately. Due to the space limit, we just list sample of most related terms for each topic in Table III.

Fig. 6 shows the topic drift in our dataset. We observe that the popular topics drifted along the time. Before 2004, the topic “Business News” was the most popular topic in email spam. After that, the most popular topic changed more frequently than before. First, the most popular topic changed to “Software Product” for around a year. And then it changed back to the topic “Business News” again. And later on, the most popular topic changes happened in the following order: “Adult Product”, “Free Product”, “Sales News”, “Free

Product”, “Newsletter”, “Official News”, “Order Information”, “Medical Product”, and “Account Information”. For each topic, it contains certain features that are attractive to certain group of users. For instance, topic “Free Product” is more attractive to users who like free stuff. Topic “Medical Product” is more attractive to users who need medical service or special medical products. Topic “Sales News” and “Order Information” are more attractive to users who like shopping. Meanwhile, as social media have interfaces with email systems normally and gain increasing popularity, email spam which have the content related to social media are growing rapidly. For instance, by investigating the content of email messages which belonging to the recent most popular topic “Account Information”, we observe that a lot of email spam have associations with social media. One example is that social media account registration email spam which contains spam URLs that camouflaged as confirmation URL links. Another

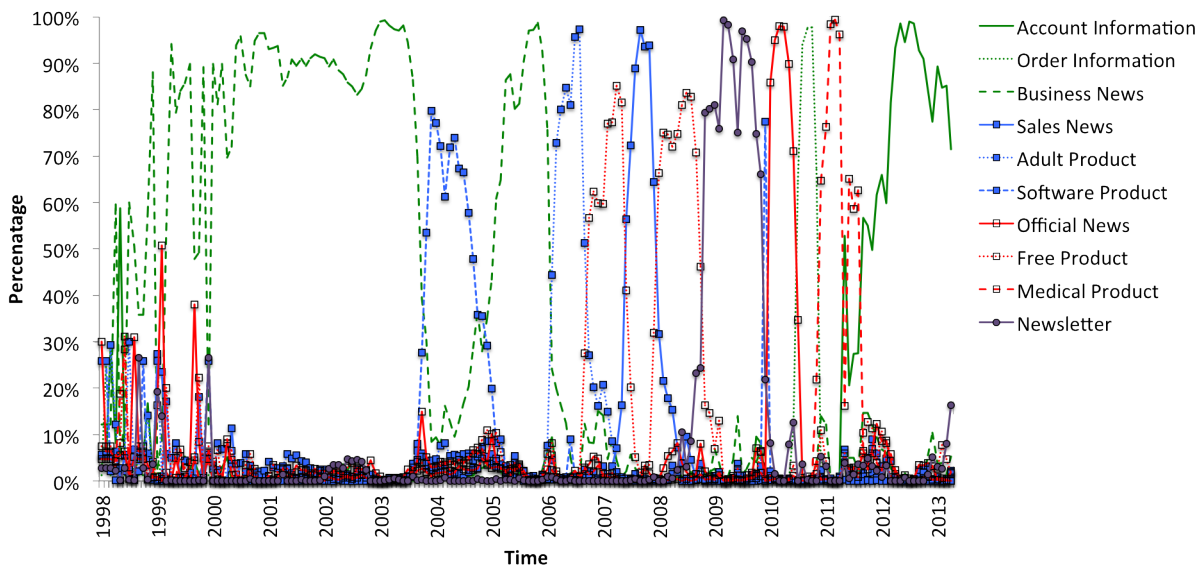


Fig. 6. Topic drift in time order (time unit: month)

TABLE II. List of topics and associated terms

Topic Name	Samples of Most Related Terms
Account Information	email important pass check account address information
Order Information	click message privacy online policy information address view order receive required
Business News	click information price free professional time link business work
Sales News	price life money make time today offer year online real world women retail deal credit people
Adult Product	world price penis back people product degree patch life make great experience enlarge
Software Product	price professional click software company copy softwares read suite online site office
Official News	united states world state national city view university government international people population
Free Product	online pills price click quality save products email item prices service offer free
Medical Product	generic save price time products medications order pharmacy home service product
Newsletter	mail click email privacy newsletter message receive view offers link receiving policy subscribed

example is social media account notifications. For example, it informs you that your account has been changed by someone and needs immediate action to reset the password, followed by the spam URL links. Thus, one possible reason why the topic “Account information” becomes popular is that a lot of spammers try to impersonate the support team of social media to steal sensitive information, such as credential and credit information, or lead users to spam or phishing web pages for further actions.

C. Network Analysis

Besides content analysis and topic modeling, we also try to find out the sending behavior changes of spammers over time through analyzing the routing network between sender and receiver. Before entering into the detail of network analysis, we will talk about data processing and some findings during the process.

For the data processing, we need to process the headers of email message to obtain the information about routing between sender and receiver. The headers which are related to the routing info are “From”, “To”, “CC”, “BCC” and “Received”. The header “From” and “To” provide the sender and receiver email addresses. The header “CC” and “BCC”

show the recipient lists in carbon copy and blind carbon copy mode. The header “Received” contains routing information from sender and receiver. First, we look into the headers “From” and “To” and intend to use them to extract the sender-to-receiver network. However, the fact is that we cannot use them in our study since most of the messages in the dataset contain forged “From” headers in one form or another, which is also mentioned in the Spam Archive dataset homepage. Although “From” header should not be trusted, we still extract top-10 domains from the “From” header to find out what are those popular domains used by spammers to set up social engineering traps for users. It is hard for users to recognize fake senders based on senders’ email address especially when the email address is belonging to the domains they trust. The list of top-10 domains is shown in Table III.

From Table III, we observe that several popular email domains are used by spammers such as “yahoo.com”, “hotmail.com”, “msn.com”, and “gmail.com”. Also some top domains are related to receiver domains such as “untroubled.org” and “dyndns.org”. It reveals that spammers were camouflaging themselves coming from the same domains as the users’ domains. In addition, some domains in the top-10 list are from countries outside US such as “163.com” which is the largest

TABLE III. List of top-10 domains

1998	1999	2000	2001	2002	2003	2004	2005
hotmail.com yahoo.com msn.com usa.net earthlink.net worldnet.att.net aol.com mailexcite.com juno.com prodigy.com	yahoo.com hotmail.com aol.com usa.net ibm.net msn.com iname.com hotmail.com bigfoot.com bigfoot.com mailcity.com	yahoo.com hotmail.com earthlink.net aol.com usa.net excite.com mail.com bigfoot.com email.com postmark.net	hotmail.com yahoo.com excite.com msn.com aol.com btamail.net.cn earthlink.net mail.com pacbell.net mail.ru	yahoo.com hotmail.com aol.com msn.com excite.com link2buy.com eudoramail.com flashmail.com netscape.net btamail.net.cn	yahoo.com hotmail.com aol.com msn.com artauction.net earthlink.net excite.com artaddiction.com juno.com artists-server.com	yahoo.com hotmail.com msn.com yahoo.co.kr aol.com attbi.com yahoo.co.jp 163.com excite.com seznam.cz netscape.net	yahoo.com hotmail.com msn.com yahoo.co.kr gmail.com yahoo.co.jp 163.com msa.hinet.net mail.com 126.com
2006	2007	2008	2009	2010	2011	2012	2013
yahoo.co.jp hotmail.com mail.ru 0451.com em.ca yahoo.com 0733.com aol.com infoseek.jp msn.com	yahoo.com dyndns.org hotmail.com yahoo.co.jp paran.com gmail.com 163.com msn.com msa.hinet.net so-net.ne.jp	dyndns.org yahoo.com adelphia.com hotmail.com gmail.com wikipedia.org earthlink.net att.net 163.com cox.net	dyndns.org homeip.net lists.untroubled.org gmail.com hotmail.com yahoo.com untroubled.org ezmlm.org em.ca mail.ru	dyndns.org yahoo.com homeip.net untroubled.org lists.untroubled.org ezmlm.org em.ca comcast.net gmail.com pfizer.com	yahoo.com dyndns.org ymail.com gmail.com mail.ru msn.com bk.ru qip.ru list.ru aol.com	yahoo.com garden.md yahoo.co.jp ageha.cc peach.6060.jp ts5558.com momoiro.cc koikoilkoi.com wakuwaku-happy.net get-c.com	yahoo.co.jp li-brooz.jp yahoo.com mixi mega.biz netstar-inc.co.uk garden.md for-dear-2013.mobi wakuwaku06.info greemmix.info docomo.ne.jp

email service domain in China. In 2013, the top domains list contains more special domains such as “.biz” which is intended for registration of domains to be used by businesses and “.mobi” which is used by mobile devices for accessing Internet resources via the Mobile Web. It indicates that spammers were spoofing the sender addresses targeting business and mobile users. Meanwhile, it proves that spammers recognize the trend of information flow in the Internet and evolve to take advantage of the trending.

Next, we investigate the header “CC” and “BCC” in email message to know whether spammers use those functions to spread email spam. The trends of “CC” and “BCC” are shown in Fig. 7.

Fig. 7 shows that spammers used more “CC” and “BCC” in the early years (1999-2004) and less in the recent years (2011-2013). One possible reason is that most spam filters have taken the number of “CC” and “BCC” as important features to detect spam. Meanwhile, people become alert to email message which contains a long recipient list in the header “CC” and “BCC” so that this type of email spam lost markets gradually.

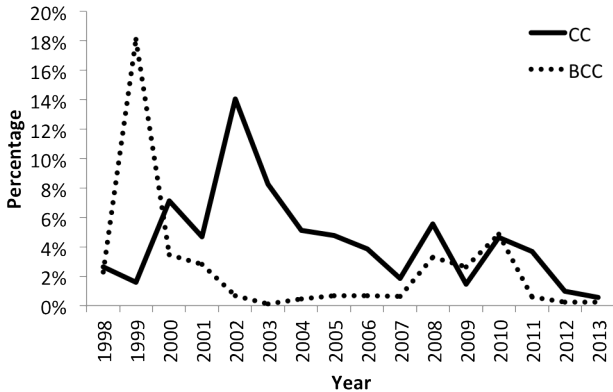


Fig. 7. Cc and Bcc trends

Based on observations above, we realize that the header “From”, “To”, “CC”, and “BCC” are not helpful in extracting routing network from email spam. To better understanding the changes in terms of spammers’ behaviors, we still need to find a way to extract the real sender and the routing information. The study of header “Received”, which is much harder to be forged, provides us the routing information such as hops’ IP addresses between sender and receiver. Therefore, we will use the header “Received” to extract sender-to-receiver IP routing information and construct routing network. The tool we used in extraction is the email module in Python [30] and the network analysis tool is the open source network visualization software Gephi [31].

During the process of extracting networks, we also collect two extra features: average hops between sender and receiver and the Geolocation distribution of sender IP addresses. The list of average hops and the Geolocation distribution of sender IP addresses over time are shown in Fig. 8 and Fig. 9 respectively.

Fig. 8 presents the trend of average hops between sender and receiver. We observe that the number of hops was increasing over time. For instance, the average hops for 1998 was only two while it became almost eight in 2013. One possible reason

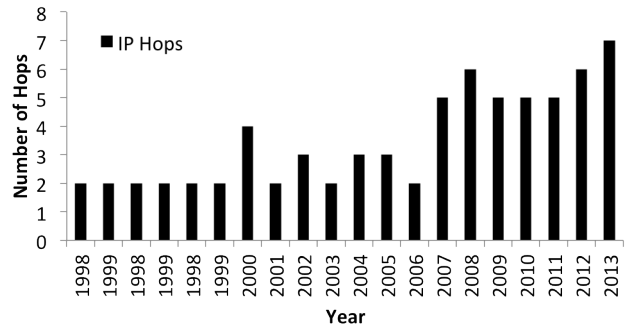


Fig. 8. Average hops between sender and receiver

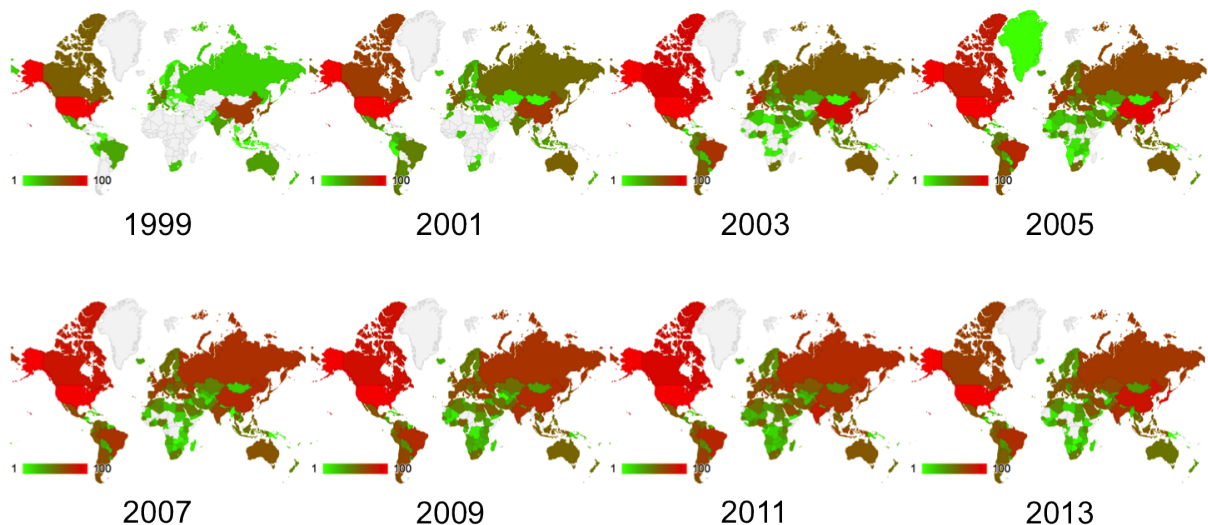


Fig. 9. Geolocation distribution of senders' IP addresses every two years (in log scale and normalized)

is that it increased the cost for spam filters to detect or trace back the senders of email spam as spammers used more hops through intermediate proxies. It also indicates that the sender-to-receiver network becomes more complicated.

Fig. 9 shows the Geolocation distribution of senders' IP addresses over time. Due to space limit, we only present the Geolocation maps every two years based on the normalized number of IP addresses coming from different countries. We use the GeoIP service provided by MaxMind [32] to do the mapping between IP address and Geolocation. Also, we employ Google Geo Chart APIs [33] to implement the map drawing. The number of IP addresses from different countries has been put into log scale and then normalized into the same range from 1 to 100. Also we use green color to label countries who had the fewest sender IP address and red color to label countries who had the more sender IP addresses. White color means that no sender IP address came from the country. Observing the maps, we have the following findings in our dataset: 1) the sender IP addresses almost come from all over the world; 2) United States has the largest number of sender IP addresses along the past fifteen years; 3) Besides United States, the distribution of sender IP addresses shows dynamic changes over time. For instance, the number of sender IP addresses coming from China kept increasing until 2007 and grew again in 2013. Also, some countries had sudden increase of sender IP addresses in particular years. For example, Canada and France had sudden increase in 2003. India had sudden increase in 2011. And Japan had sudden increase in 2013. It indicates that spammers used global email service servers and also kept changing the traffic from different countries.

Next, we extract the networks from our dataset for each year and use three major metrics to measure the complexity of them. The three metrics are network diameter (the longest of all the calculated shortest paths in a network), average degree (average number of edges connected to or from one node), and average clustering coefficient (a measure of degree to which nodes in a network tend to cluster together). The result of measurement is shown in Fig. 10. Since the data for 1998 and

2013 does not cover the whole year, we only list the results from 1999 to 2012.

Fig. 10 shows the three metrics comparison from 1999 to 2012. The values of them have the increasing trend overall but fluctuations existed along the time. Network diameter became more stable after 2007 and it is the same for the metrics average degree and average clustering coefficient. Those metrics kept staying at high value in terms of complexity of network.

For the purpose of better visualization, we remove those nodes whose degree is lower than certain threshold. And also due to the space limit, we only present the network graph every five years (1998, 2003, 2008, and 2013) in Fig. 11. For 1998 and 2003, we keep the nodes whose degree is greater than 3. While for 2008 and 2013, we keep the nodes whose degree is greater than 10. The reason is that too many node overlaps occur if we choose the threshold 3 for 2008 and 2013.

Fig. 11 shows the sender-to-receiver routing network based on the IP addresses extracted from email header "Received". We observe that the complexity of graph increases explicitly along the time. For 2013, even that it only contains half year data, the routing network has already shown much more complicated than the routing network in 2008.

V. DISCUSSION

Our large-scale evolutionary study on email spam dataset in a long period of time shows the trend of email spam business. Although the volume of email spam had a slight drop in recent years, we cannot conclude that email spam business is dying and email spam filters have won the battle against spammers. Through intensive analysis including content analysis, topic modeling and network analysis, we demonstrated that the battle is still ongoing and even worse since spammers become more sophisticated and capricious. Moreover, our study still have the following limitations and future work to do.

The dataset we used does not cover all the email spam over the fifteen years, which may influence the accuracy of our results, especially for the portion in the early years such as 1998-2000 that contains small number of email spam. Also,

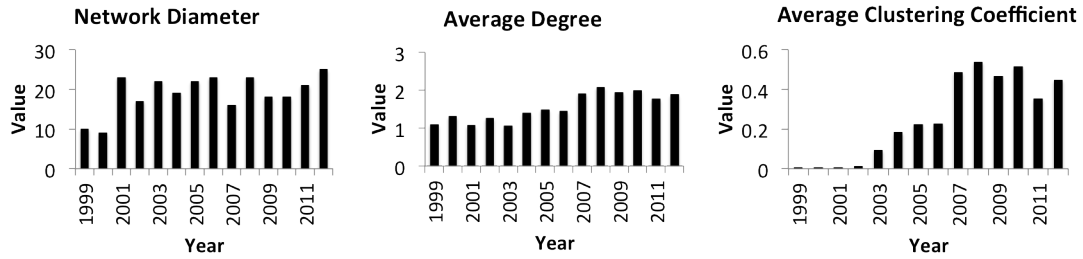


Fig. 10. The comparison of three metrics from 1999 to 2012

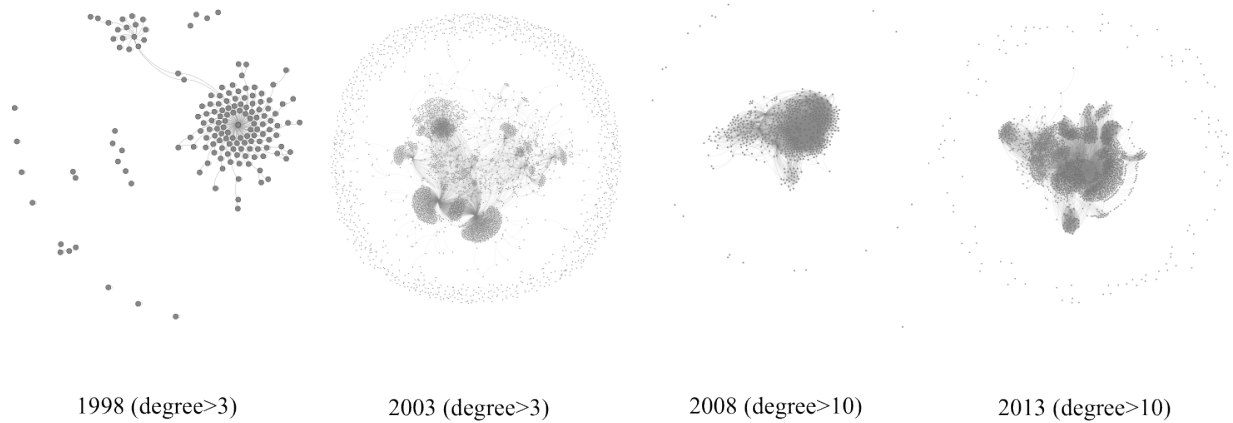


Fig. 11. Sender-to-receiver routing networks every five years from 1998 to 2013

the bait email addresses used in data collection may cause some biases in the dataset. For example, the domain of the email address may result in that spammers forge their email addresses to the same domain.

Besides the limitation on dataset, we also have limitation on our analysis. In the topic modeling analysis, we set up the number of topics to 10 that may influence the result of topic modeling. If we change the number of topics to larger value, the result may be more accurate and fine-grained. But it should not conflict with our conclusion that the topic drift occurs frequently over time. We will take the fine-grained analysis as future work. Additionally, in the network analysis, we used the study of the header “Received” to extract sender-to-receiver network. But we cannot guarantee that no forged information exists in the header “Received”. Spammers also have some techniques to spoof the header “Received” but the portion of forged headers is low since it costs spammers a lot and has certain strict requirements to meet. We will also look into the further validation work in the future.

VI. RELATED WORK

Our work mainly involves three lines of research work: email spam detection, analysis approach on email data, and evolutionary study of spam.

Email spam detection has been studied by lots of researchers in different directions. For instance, Carreras et al. [34] applied boosting trees to filter out email spam. Wang et al. [35]

used heuristic feature selection techniques to improve the performance of email spam filtering. Chan et al. [36] co-trained with a single natural feature set in email classification. Liu et al. [37] adopted multi-field learning for email spam classification. Sculley et al. [38] used relaxed online SVMs for email spam filtering. Besides those machine learning techniques, more researchers tried other kinds of detection methods. Attenberg et al. [39] introduced collaborative email spam filtering with the hashing trick. Balakumar et al. [40] offered ontology based classification of email. Dasgupta et al. [41] combined similarity graphs to enhance email spam filtering. Jung et al. [42] used DNS black lists and spam traffic to detect email spam. Ramachandran et al. [43] filtered email spam with behavioral blacklisting. Clayton et al. applied extrusion detection in stopping email spam by observing distinctive email traffic patterns. Xie et al. [44] provided an effective defense approach against email spam laundering. Additionally, researchers also have used email spam to help detecting other types of spam. For instance, Zhuang et al. [45] developed an approach to map botnet membership using traces of spam email. Webb et al. [46] identified an interesting link between email spam and Web spam and used it to extract large Web spam samples from the Web. Wang et al. [47] demonstrated the relationship among different formats of social spam including user profile spam, message spam and Web spam, in which message spam contain email spam.

Analysis approach on email data is another focus of researchers. Bird et al. [48] constructed social networks of

email correspondents to address some interesting questions such as the social status of different types of participants and the relationship of email activity and other activities. McCallum et al. [49] illustrated experimental study on Enron and academic email to discover topic and role in social networks from emails, in which the model builds on Latent Dirichlet Allocation (LDA) and the Author-Topic (AT) model. Culotta et al. [50] presented an end-to-end system that extracts a user's social network and its members contact information given the user's email inbox.

Research work on evolutionary study of spam is close to this paper. Pu et al. [15] presented a study on dataset collected from Spam Archive and focused on two evolutionary trends: extinction and existence. Irani et al. [17] studied the evolution of phishing email messages and classified them into two groups: flash attacks and non-flash attacks. Wang et al. [18] compared two large Web spam corpus: Webb spam corpus 2006 and Webb spam corpus 2011 and shown the trending of Web spam. Chung et al. [51] and Fetterly et al. [52] also have done intensive study on evolution of web spam. Guerra et al. [19] investigated how the popularity of spam construction techniques changes when filters start to detect them and determined automatically techniques that seemed more resistant than others.

The evolution of spamming techniques shows the increasing sophistication of spammers. Our work focuses on tactics changes of email spam over time and inspires more researchers to work on email spam detection collaboratively.

VII. CONCLUSIONS

We introduced a long-term evolutionary study on large scale email spam corpus – Spam Archive dataset, which contains over 5 million email messages from 1998 to 2013. Besides content analysis of email spam including n -grams analysis, we adopted topic modeling and network analysis techniques to investigate topic drift and increasing complexity of sending behaviors of spammers.

In the topic modeling, we clustered our dataset based on LDA model and categorized them into ten topics: “Account Information”, “Order Information”, “Business News”, “Sales News”, “Adult Product”, “Software Product”, “Official News”, “Free Product”, “Medical Product”, and “Newsletter” based on the most related terms associated. The result shows spammers changed topics over time and also made the topics attractive to users. In the network analysis, we presented social engineering attacks from spammers by observing senders' domains. After studying the header “Received”, we extracted sender IP addresses and the sender-to-receiver routing networks from the dataset. The Geolocation distribution of senders' IP addresses shows that spammers employed the servers all over the world and dynamically switched locations among different countries. Moreover, we chose three metrics: network diameter, average degree, and average clustering coefficient to measure the complexity of routing networks, showing that the sending behaviors of spammers are becoming more complicated and harder to track.

To sum up, email spam business is becoming more sophisticated along the time and the spammers behind it evolve into

more capricious in the ongoing battle with spam filters.

ACKNOWLEDGEMENTS

This research has been partially funded by National Science Foundation by CNS/SAVI (1250260), IUCRC/FRP (1127904), CISE/CNS (1138666), RAPID (1138666), CISE/CRI (0855180), NetSE (0905493) programs, and gifts, grants, or contracts from DARPA/I2O, Singapore Government, Fujitsu Labs, and Georgia Tech Foundation through the John P. Imlay, Jr. Chair endowment. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation or other funding agencies and companies mentioned above.

REFERENCES

- [1] S. Whittaker, V. Bellotti, and J. Gwizdka, “Email in personal information management,” *ACM Communications*, vol. 49, no. 1, pp. 68–73, Jan. 2006.
- [2] R. Clayton, “Email traffic: a quantitative snapshot,” in *the 4th Conference on Email and Anti-Spam (CEAS 2007)*, Mountain View, CA, USA, July 2007.
- [3] L. F. Cranor and B. A. LaMacchia, “Spam!” *ACM Communications*, vol. 41, no. 8, pp. 74–83, Aug. 1998.
- [4] MAAWG, “Email Metrics Report 2011,” Tech. Rep., November 2011. [Online]. Available: http://www.maawg.org/sites/maawg/files/news/MAAWG_2011_Q1Q2Q3_Metrics_Report_15.pdf
- [5] J. Goodman, G. V. Cormack, and D. Heckerman, “Spam and the ongoing battle for the inbox,” *ACM Communications*, vol. 50, no. 2, pp. 24–33, Feb. 2007.
- [6] M. Sahami, S. Dumais, D. Heckerman, and E. Horvitz, “A bayesian approach to filtering junk e-mail,” in *Learning for text categorization: papers from the 1998 workshop*, 1998.
- [7] A. Cournane and R. Hunt, “An analysis of the tools used for the generation and prevention of spam,” *Computers & Security*, vol. 23, no. 2, pp. 154 – 166, 2004.
- [8] Z. Gyongyi and H. Garcia-Molina, “Web spam taxonomy,” Stanford InfoLab, Technical Report 2004-25, March 2004.
- [9] S. Y. Park, J.-T. Kim, and S.-G. Kang, “Analysis of applicability of traditional spam regulations to voip spam,” in *Advanced Communication Technology, 2006. ICACT 2006. The 8th International Conference*, vol. 2, 2006, pp. 3 pp.–1217.
- [10] P. Hayati, V. Potdar, A. Talevski, N. Firoozeh, S. Sarenche, and E. Yeganeh, “Definition of spam 2.0: New spamming boom,” in *Proceedings of the 4th IEEE International Conference on Digital Ecosystems and Technologies (DEST)*, 2010, pp. 580–584.
- [11] D. Fallows, “Spam. how it is hurting email and degrading life on the internet,” the Pew Internet & American Life project, Washington, DC, USA, Technical Report, October 2003.
- [12] N. Dalvi, P. Domingos, Mausam, S. Sanghai, and D. Verma, “Adversarial classification,” in *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, ser. KDD '04. New York, NY, USA: ACM, 2004, pp. 99–108.
- [13] D. Chinavle, P. Kolari, T. Oates, and T. Finin, “Ensembles in adversarial classification for spam,” in *Proceedings of the 18th ACM conference on Information and knowledge management*, ser. CIKM '09. New York, NY, USA: ACM, 2009, pp. 2015–2018.
- [14] B. Biggio, G. Fumera, and F. Roli, “Evade hard multiple classifier systems,” in *Applications of Supervised and Unsupervised Ensemble Methods*, ser. Studies in Computational Intelligence, O. Okun and G. Valentini, Eds. Springer Berlin Heidelberg, 2009, vol. 245, pp. 15–38.
- [15] C. Pu and S. Webb, “Observed trends in spam construction techniques: A case study of spam evolution,” in *Proceedings of the Third Conference on Email and Anti-Spam (CEAS 2006)*, Mountain View, CA, USA, July 2006.
- [16] T. Fawcett, ““in vivo” spam filtering: a challenge problem for kdd,” *SIGKDD Explor. Newsl.*, vol. 5, no. 2, pp. 140–148, Dec. 2003.
- [17] D. Irani, S. Webb, J. Giffin, and C. Pu, “Evolutionary study of phishing,” *eCrime Researchers Summit*, 2008, pp. 1–10, 2008.

- [18] D. Wang, D. Irani, and C. Pu, "Evolutionary study of web spam: Webb spam corpus 2011 versus webb spam corpus 2006," in *Proceedings of the 8th International Conference on Collaborative Computing: Networking, Applications and Work-sharing (CollaborateCom)*, Pittsburgh, PA, USA, October 2012, pp. 40–49.
- [19] P. Guerra and D. Guedes, "Exploring the spam arms race to characterize spam evolution," in *Proceedings of the 8th Annual Collaboration, Electronic messaging, Anti-Abuse and Spam Conference (CEAS 2010)*, Redmond, Washington USA, July 2010.
- [20] "Kaspersky lab," 2013. [Online]. Available: <http://usa.kaspersky.com/>
- [21] K. Rapoza, "The dying business of email spam," <http://usa.kaspersky.com/about-us/press-center/in-the-news/dying-business-email-spam>, 2012.
- [22] D. Gudkova, "Kaspersky security bulletin: Spam evolution 2012," http://www.securelist.com/en/analysis/204792276/Kaspersky_Security_Bulletin_Spam_Evolution_2012, 2012.
- [23] "Untroubled website," 2013. [Online]. Available: <http://untroubled.org/spam/>
- [24] "Electronic communications and transactions act, 2002," 2002. [Online]. Available: http://www.internet.org.za/ect_act.html
- [25] "Multipurpose internet mail extensions (mime) part one: Format of internet message bodies," 1996. [Online]. Available: <http://tools.ietf.org/html/rfc2045>
- [26] "Text:ngrams - flexible ngram analysis (for characters, words, and more)," 2013. [Online]. Available: <http://search.cpan.org/dist/Text-Ngrams/Ngrams.pm>
- [27] "Defensive aids sub system (dass)," 2013. [Online]. Available: <http://www.eurofighter.com/capabilities/technology/sensor-fusion/defensive-aids-sub-system.html>
- [28] A. K. McCallum, "Mallet: A machine learning for language toolkit," 2002. [Online]. Available: <http://mallet.cs.umass.edu>
- [29] D. M. Blei, A. Ng, and M. Jordan, "Latent dirichlet allocation," *JMLR*, vol. 3, pp. 993–1022, 2003.
- [30] "Python: email an email and mime handling package," 2013. [Online]. Available: <http://docs.python.org/2/library/email>
- [31] "Gephi, an open source graph visualization and manipulation software," 2013. [Online]. Available: <http://gephi.org>
- [32] "Maxmind - ip geolocation and online fraud prevention," 2013. [Online]. Available: <http://www.maxmind.com/en/home>
- [33] "Visualization: Geochart - google charts google developers," 2013. [Online]. Available: <https://developers.google.com/chart/interactive/docs/gallery/geochart>
- [34] X. Carreras and L. Marquez, "Boosting trees for anti-spam email filtering," *arXiv preprint cs/0109015*, 2001.
- [35] R. Wang, A. Youssef, and A. Elhakeem, "On Improving the Performance of Spam Filters Using Heuristic Feature Selection Techniques," in *Proceedings of 23rd Biennial Symposium on Communications, 2006*. Ieee, 2006, pp. 227–230.
- [36] J. Chan, I. Koprinska, and J. Poon, "Co-training with a Single Natural Feature Set Applied to Email Classification," in *Proceedings of IEEE/WIC/ACM International Conference on Web Intelligence (WI'04)*. Ieee, 2004, pp. 586–589.
- [37] W. Liu and T. Wang, "Multi-field learning for email spam filtering," in *Proceeding of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. New York, NY, USA: ACM Press, 2010, p. 745.
- [38] D. Sculley and G. Wachman, "Relaxed online SVMs for spam filtering," in *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, April 2007, pp. 415–422.
- [39] J. Attenberg, K. Weinberger, and A. Dasgupta, "Collaborative Email-Spam Filtering with the Hashing Trick," in *CEAS*, 2009, pp. 1–4.
- [40] M. Balakumar and V. Vaidehi, "Ontology based classification and categorization of email," in *Proceedings of Signal Processing, Communications and Networking*, 2008, pp. 199–202.
- [41] A. Dasgupta, M. Gurevich, and K. Punera, "Enhanced email spam filtering through combining similarity graphs," in *Proceedings of the fourth ACM international conference on Web search and data mining*. New York, NY, USA: ACM Press, 2011, p. 785.
- [42] J. Jung and E. Sit, "An empirical study of spam traffic and the use of DNS black lists," in *Proceedings of the 4th ACM SIGCOMM conference on Internet measurement*. New York, NY, USA: ACM Press, 2004, p. 370.
- [43] A. Ramachandran, N. Feamster, and S. Vempala, "Filtering spam with behavioral blacklisting," in *Proceedings of the 14th ACM conference on Computer and communications security*. New York, NY, USA: ACM Press, 2007, p. 342.
- [44] M. Xie, H. Yin, and H. Wang, "An effective defense against email spam laundering," in *Proceedings of the 13th ACM conference on Computer and communications security*. New York, NY, USA: ACM Press, 2006, p. 179.
- [45] L. Zhuang, J. Dunagan, D. Simon, and H. Wang, "Characterizing Botnets from Email Spam Records," in *Proceedings of the first USENIX workshop on large-scale exploits and emergent threats (LEET 08)*, 2008.
- [46] S. Webb, J. Caverlee, and C. Pu, "Introducing the webb spam corpus: Using email spam to identify web spam automatically," in *Proceedings of the Third Conference on Email and Anti-Spam (CEAS 2006)*, Mountain View, CA, USA, July 2006.
- [47] D. Wang, D. Irani, and C. Pu, "A social-spam detection framework," in *Proceedings of the 8th Annual Collaboration, Electronic messaging, Anti-Abuse and Spam Conference (CEAS 2011)*, Perth, Australia, September 2011, pp. 46–54.
- [48] C. Bird, A. Gourley, and P. Devanbu, "Mining email social networks," in *the 2006 international workshop on Mining software repositories*, 2006, pp. 137–143.
- [49] A. McCallum, X. Wang, and A. Corrada-Emmanuel, "Topic and role discovery in social networks with experiments on enron and academic email," *J. Artif. Intell. Res.(JAIR)*, vol. 30, pp. 249–272, 2007.
- [50] A. Culotta, R. Bekkerman, and A. McCallum, "Extracting social networks and contact information from email and the web," in *Proceedings of the First Conference on Email and Anti-Spam (CEAS 2004)*, 2004.
- [51] Y. Chung, "A study on the evolution and emergence of web spam," Ph.D. dissertation, Univ. of Tokyo, Tokyo, Japan, 2011.
- [52] D. Fetterly, M. Manasse, M. Najork, and J. Wiener, "A large-scale study of the evolution of web pages," in *Proceedings of the 12th international conference on World Wide Web*, ser. WWW '03, New York, NY, USA, 2003, pp. 669–678.