# Hybrid Spam Filtering for Mobile Communication

Ji Won Yoon
Robotics Research Group
Engineering Department,
University of Oxford, UK
jwyoon@robots.ox.ac.uk

Hyoungshick Kim
Security Group
Computing Laboratory,
University of Cambridge, UK
hk331@cl.cam.ac.uk

Jun Ho Huh
Software Engineering Group
Computing Laboratory,
University of Oxford, UK
jun.ho.huh@comlab.ox.ac.uk

## ABSTRACT

Spam messages are an increasing threat to mobile communication. Several mitigation techniques have been proposed, including white and black listing, challenge-response and content-based filtering. However, none are perfect and it makes sense to use a combination rather than just one. We propose an anti-spam framework based on the hybrid of content-based filtering and challenge-response. There is the trade-offs between *accuracy* of anti-spam classifiers and the *communication overhead*. Experimental results show how, depending on the proportion of spam messages, different filtering parameters should be set.

## Keywords

Anti-spam filtering, a threshold sensitivity problem, Uncertain area in decision, Human Interaction

## 1. INTRODUCTION

Short Message Service (SMS) and Multimedia Messaging Service (MMS) are a popular means of mobile communication. Texting costs have decreased continuously over the years (to an extent of free texting) whereas the bandwidth for communication has increased dramatically. Such trends have attracted a large number phishing and spamming attacks using text messages. In particular, spam messages containing pornographic or promotive materials are an emerging phenomenon, and they have caused a significant level of inconvenience for users. These are now prevalent in Korea, Japan and China and prone to spread across countries where mobile communication is popular. Statistics for 2008 show that a user in China, on average, receives 8.29 SMS spam per week [1].

Much of the existing research into anti-spam solutions, however, has focused on the protection of emails in the context of the Internet. Some of the popular methods include white and black listing, digital signature, postage control, address management, collaborative and content-based filtering [2, 3, 4, 5, 6]. Different characteristics between emails and text messages make it harder for one to apply such approaches directly in mobile networks and analyze the results. Each approach has its own set of drawbacks and does not improve much if used alone. For example, the extra traffic required to perform challenge-response needs to be minimized (or needs to be compensated for) as it is more expensive to use the bandwidth in mobile networks. This issue is capable being addressed by content-based filtering: obvious spam would be filtered first to reduce the number of messages subject to challenge-response. In this paper, we propose an anti-spam framework based on the combination of these two methods. We attempt to reduce a great number of high-volume spamming, and as a result, minimize the extra amount of bandwidth that would be required. Given a reasonable filtering algorithm, we show that, ultimately, less bandwidth (than freely allowing high-volume spam) will be used with our method.

The remainder of the paper is organized as follows. In Section 3, we describe a hybrid spam filtering framework. Section 4 evaluates the performance of our hybrid method based on two measures, the traffic usage and the accuracy. Finally, in Sections 5 we discuss the contribution of this paper and the remaining work.

## 2. EXISTING SOLUTIONS

Content-based filtering solutions have been proved to be effective against emails, which are typically larger in size compared to text messages. Abbreviations and acronyms are used more frequently in SMS and they increase the level of ambiguity. [7] propose binary classification and filtering methods for short messages in a Bayesian scheme. Classification rules are defined and extended from general patterns identified in past spam. However, adaptive schemes as such are weak against innovative attacks where strategies constantly evolve to manipulate classification rules. Filtering alone will not be sufficient to detect spam.

Many anti-spam solutions [1, 8] have been suggested based on a challenge-response protocol. A message sender needs to verify that they are a legitimate sender by answering the challenge message (e.g. through a web interface) before their message is forwarded to the recipient. The sender authenticates themselves as a human-user by answering a simple turing test for which a machine cannot easily understand. Nevertheless, the protocol has often been criticized for extra user interaction and traffic used. There might also be a significant overhead in storing and managing challenge messages.

Our goal is to develop a solution that ultimately mini-

mizes the usage of network bandwidth by discouraging high-volume spamming. We believe the extra traffic required to perform challenge-response can be compensated if a large amount of spamming attacks can be reduced as a result. In our approach the challenge-response protocol classifies machine-generated spam. We also use the filtering method to reduce the number of messages that need to be verified. Simulation results in later sections show that this hybrid approach is capable of controlling high-volume spam and the traffic usage.

## 3. A HYBRID FRAMEWORK

Text messages are classified into three different regions using a content-based filtering method: normal, uncertain and spam. A filtering method cannot deal with *uncertain* messages; therefore, we use a challenge-response protocol to classify the uncertain messages into normal and spam regions (see Fig. 1). We assume that the majority of spam messages are generated by machines.

A human verification mechanism (in the form of challenge-response) is added to a common filtering scheme to detect whether an uncertain message falls into the normal or the spam region. A message center (owned by an operator) sends a challenge query to check if the sender is an individual or a machine. The sender responds by answering the query and the center compares the returned value against the known correct value. If the values match, the message is classified as *normal*, else, as *spam*. We are interested in further classification of this uncertain region. A message
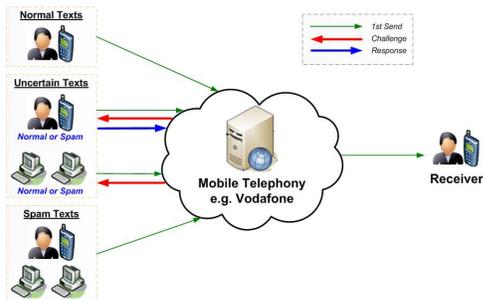


**Figure 1: Hybrid Spam Filtering Overview**

center is given the full responsibility of running the framework due to the following reasons. Firstly, it should reduce the traffic usage by filtering spam as early as possible, before forwarding them to the recipient. Secondly, using the challenge-response protocol, the center will be able to collect an enormous amount of sample data in real time; these can be used to develop highly effective classifiers and continuously improve the performance of filtering algorithms. Lastly, it would be difficult to install and maintain a homogeneous anti-spam software on all mobile devices; instead we rely on one solution deployed in a message center.

### 3.1 Uncertain Region

If we assume there are only two regions, normal and spam, a filtering method will use binary classification. Suppose that we have a probabilistic model for the anti-spam classifier as a posterior distribution $Pr(c = \text{normal}|y)$. This is the probability that a message is normal: $c$ and $y$ denote realization of random variables for a class and a message,



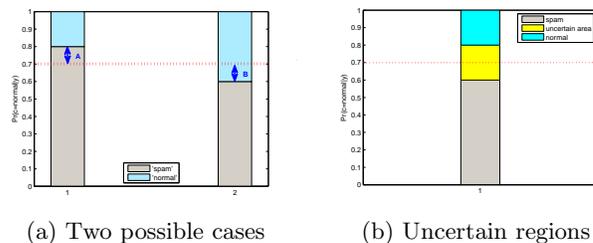(a) Two possible cases          (b) Uncertain regions

**Figure 2: (a) Two possible cases: $h > \tilde{h}$ (case 1) and $h < \tilde{h}$ (case 2) for a given ground truth $\tilde{h}$ (red dot line) and (b) Modified classification embedding uncertain area given a ground truth $\tilde{h}$ (red dot line)**

respectively. The odd ratio of the posterior is used to obtain a measurable classification by $O_{post} = \frac{Pr(c=\text{normal}|y)}{Pr(c=\text{spam}|y)}$. If $O_{post} > 1$, a message is classified as normal; otherwise, as spam. Alternatively, we can simply use a threshold based approach in the posterior distribution. If $Pr(c = normal|y)$ is close to one, a message is likely to be normal; if close to zero, it is likely to be spam. Let $\bar{c} = f(y, h)$ be a spam filter where $\bar{c}$ and $h$ are an output and a given threshold, respectively. This filter would work with the following rules:

$$\bar{c} = f(y,h) = \left\{ \begin{array}{ll} \text{normal} & \text{if } Pr(c = \text{normal}|y) \geq h \\ \text{spam} & \text{if } Pr(c = \text{normal}|y) < h \end{array} \right. \quad (1)$$

This separates normal messages from spam (odd ratio approach is a special case where $h = 0.5$). The main problem with this approach is finding a proper threshold; because the threshold for ground truth $\tilde{h}$ is unknown, there are two possible cases as shown in Fig. 2. If $h$ is higher than $\tilde{h}$, some of the normal messages in region $A$ may be classified as spam. If $h$ is lower than $\tilde{h}$, some of the spam in region $B$ may bypass the filter and reach the recipients. Such a threshold problem will always be present in classification: it is almost impossible to find the underlying $\tilde{h}$; and anti-spam software companies are likely to use strategies based on their own experiences. The sensitivity problem can be resolved by introducing an uncertain region with two thresholds (see Fig. 2-b). These can be implemented as the upper and lower boundaries of a traditional threshold system. There are three labels: spam, uncertain area, and normal; we focus on the uncertain area. Spam and normal regions are classified as in the traditional system. Only the messages that fall into the uncertain area are checked using challenge-response. In the next section we describe the protocols in detail.

### 3.2 Challenge-Response Protocols

We assume that there is a turing test available with a low probability of producing false positives and negatives. Completely Automatic Public Test to tell Computer and Humans Apart (CAPTCHA) is a commonly used one: it generates pattern matching problems for which a human can easily recognize and a machine cannot. An automated program that generates thousands of spam messages will not be capable of answering a CAPTCHA based challenge, which might be a graphical image containing a faint typeface. If the response is correct, there is a high probability that the sender is a human. A CAPTCHA can be designed in a flexible manner using different media forms such as an image, an audio file and a text [9]. Their implementation details are beyond the

scope of this paper.

A number of challenge-response protocols have already been proposed [1, 8]. However, these focus only on the implementation issues without considering the security model and the cryptographic details. This section defines our security models and describes a number of possible protocols in line with them. There are several issues we need to consider before designing the protocols. Firstly, in dealing with spam, message authentication and integrity are important; whereas, confidentiality is not. Secondly, text messages are usually unencrypted and unsigned; it is possible to tamper with them during transmission. Thirdly, security properties of the communication channel between a message center and a sender need to be defined; this channel might or might not be an authenticated one. Lastly, management of the session information between all trusted pairs while performing challenge-response, would impose a huge storage overhead on a message center; there might be more than one message center sharing this information; and it might or might not be stored in the center. Mindful of these security and scalability issues, we propose four different protocols: protocols 3 and 4 have been designed with the assumption of an authenticated channel, and protocols 1 and 2 have not; moreover, protocols 1 and 3 assume that a message center manages the session information, and the others do not.

### 3.2.1 Notations

The symbols $S$ and $R$ represent a sender and a recipient, respectively. $M$ represents a mobile message center, $T$ a timestamp, $N$ a nonce, $K$ a key and $K^{-1}$ its inverse. In a symmetric crypto-system such as AES, $K$ and $K^{-1}$ are always equal. We use $\{P\}_K$ for a plain text message $P$ encrypted with $K$. $H$ is a one-way hash function. The subscript $m$ in $K_m$ implies that $K_m$ is $M$'s public key. In addition, $ms$ in $K_{ms}$ shows that $K_{ms}$ is intended for communication between $M$ and $S$.

A sender's ability to respond to a challenge depends on knowing and interpreting a key, $K_c^{-1}$. A non-authorized sender (e.g. a program sending spam) will not be able to interpret and gain information about $K_c^{-1}$; this key serves to identify machine-generated spam. For simplicity, encryption algorithms are not considered in our protocols.

### 3.2.2 Protocols

In protocol 1, the message center, $M$, maintains the session information.

[**Protocol 1**]
(M1) $S \longrightarrow M :$  $S, R, P$
(M2) $M \longrightarrow S :$  $M, S, \{K_{ms}\}_{K_c}, \{H(S, R, P), N\}_{K_{ms}}$
(M3) $S \longrightarrow M :$  $S, M, \{H(S, R, P), N+1\}_{K_{ms}}$

Before sending message 1, $S$ stores $R$ and $P$ to prevent message modification attacks. After receiving message 1, $M$ generates $K_{ms}$ and stores $(S, R, P, K_{ms}, N)$ as the session information. $K_{ms}$ is protected with $K_c$. An image CAPTCHA would be one way of protecting $K_{ms}$ against spam programs. After receiving message 2, $S$ decrypts $\{K_{ms}\}_{K_c}$ by answering the challenge (their ability to interpret $K_c^{-1}$). $S$ then decrypts $H(S, R, P)$ and $N$ using $K_{ms}$. $S$ compares $H(S, R, P)$ against the previously stored values. $S$ terminates the protocol if these values do not match; otherwise, $S$ generates $\{H(S, R, P), N+1\}_{K_{ms}}$ by $K_{ms}$ and sends it to $M$. After receiving message 3, $M$ verifies $\{H(S, R, P), N+1\}_{K_{ms}}$. If it is valid, $M$ forwards the stored message $(S, R, P)$

to $R$. Finally, $M$ deletes the session information.

Users will be frustrated if challenge-response happens too often. We use a timestamp, $T$, to solve this problem. After receiving message 3, $M$ maintains a session information $(S, R, P, K_{ms}, T)$ between $S$ and $R$ for a given time interval. $M$ checks the validity of $K_{ms}$ using the session information and a policy that defines the lifetime of $K_{ms}$.

The main drawback of this protocol is that $M$ has to bear the huge overhead of maintaining the session information. We describe another protocol which solves this issue by using authorized tokens instead:

[**Protocol 2**]
(M1) $S \longrightarrow M :$  $S, R, P$
(M2) $M \longrightarrow S :$  $M, S, \{K_{ms}\}_{K_c}, \{H(S, R, P)\}_{K_{ms}},$
  $\{K_{ms}, H(S, R), T\}_{K_m^{-1}}$
(M3) $S \longrightarrow M :$  $S, R, \{P\}_{K_{ms}}, \{K_{ms}, H(S, R), T\}_{K_m^{-1}}$

The key difference is the use of $\{K_{ms}, H(S, R), T\}_{K_m^{-1}}$ (which can only be generated by $M$) as the authorization token for verifying a response. $M$ checks whether $S$ is authorized by looking at $\{K_{ms}, H(S, R), T\}_{K_m^{-1}}$. Using this token, $S$ can just send message 3 alone, including a new text ($P'$), within the lifetime of $T$:

(M1) $S \longrightarrow M :$  $S, R, \{P'\}_{K_{ms}}, \{K_{ms}, H(S, R), T\}_{K_m^{-1}}$

In these protocols, however, $S$ cannot find out where the challenge comes from. In an attempt to solve this problem, we assume there is an authenticated channel between $M$ to $S$, and $M$'s public key $K_m$ is securely installed in a mobile device owned by $S$; perhaps during the process of manufacturing. We describe the following two protocols based on these assumptions:

[**Protocol 3**]
(M1) $S \longrightarrow M :$  $S, R, P$
(M2) $M \longrightarrow S :$  $M, S, \{\{K_{ms}\}_{K_c}, N\}_{K_m^{-1}}$
(M3) $S \longrightarrow M :$  $S, R, \{N+1\}_{K_{ms}}$

In protocol 3, $M$ maintains the session information, $(S, R, P, K_{ms}, N)$. When message 2 arrives, $S$ verifies the signature on $\{\{K_{ms}\}_{K_c}, N\}_{K_m^{-1}}$. $S$ does not respond if the signature is unknown.

[**Protocol 4**]
(M1) $S \longrightarrow M :$  $S, R, P$
(M2) $M \longrightarrow S :$  $M, S, \{\{K_{ms}\}_{K_c}, H(S, R), T\}_{K_m^{-1}},$
  $\{P\}_{K_m^{-1}}$
(M3) $S \longrightarrow M :$  $S, R, \{P\}_{K_{ms}},$
  $\{\{K_{ms}\}_{K_C}, H(S, R), T\}_{K_m^{-1}}$

Protocol 4 uses $\{\{K_{ms}\}_{K_C}, H(S, R), T\}_{K_m^{-1}}$ as the authorized token. Our protocols are likely to be compatible with existing devices since the majority already have built-in encryption and hash functions.

## 3.3 Observations

### 3.3.1 Upgrading Protocols

A message is always sent to the message center of the contracted operator first. If the message is directed at someone contracted to a different operator, it is forwarded to another message center before reaching the receiver's handset [10]. This means if one of the message centers decides not to use

our framework, all uncertain texts delivered via that center would bypass the spam filter. It would be the weakest point (and the only route needed) for an attack. Hence, all existing message centers would have to support the new protocol. While this is a large change and a challenging one, operator-sponsored forums like OMTP (Open Mobile Terminal Platform), are working with key mobile operators to unify and recommend mobile terminal requirements [11]. With the increasing number of spam texts, it seems likely that the ability to filter machine-generated uncertain texts will persuade operators into upgrading their systems.

### 3.3.2 Performance

If there are too many messages subject to challenge-response, its overhead will dominate; for example, sending an image CAPTCHA is a huge overhead to authenticate a 100 character text message. Future work may look at adding a 'bypass' to the hybrid: for example, if a message begins with a user-settable password (typically the recipient's name, but changeable), then it should be automatically treated as normal. As the uncertain region becomes smaller, we expect the performance of our framework to improve.

### 3.3.3 Usability Issues

Adapting CAPTCHA methods will have implications on the usability. A device might not have the capability to display an image CAPTCHA to a readable standard; also a mobile user might find it difficult to verify an audio CAPTCHA due to the background noise. These issues however, are likely to be resolved with technological advances. According to the Kelsey Group's second annual study on mobile use, more users than ever own internet-enabled smartphones: phones which provide advanced information accessing functions.

## 4. EVALUATION

### 4.1 Description of Datasets

In order to measure the performance of our framework, we have generated synthetic datasets. Suppose that there are $N$ number of sent messages (we set $N = 5000$). We will use $p$ and $q$ to show the normal to spam proportion where $p + q = 1$, and $p$ and $q$ are non-negative numbers (in reality, there would be many operators with different proportions). Let $\kappa$ be a random variable generated from an existing filtering method: $\kappa = Pr(c = normal|y)$. For an artificial dataset, we build a mixture model given by

$$
\begin{aligned}
p(\kappa|\lambda) &= p(\kappa|c = normal, \lambda)p(c = normal|\lambda) \\
&+ p(\kappa|c = spam, \lambda)p(c = spam|\lambda) \quad (2)
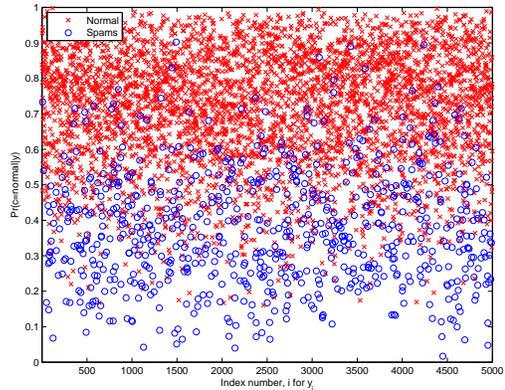\end{aligned}
$$

where $\lambda$ denotes a set of hyper-parameters which control parameters. We assume $c_i$ is generated from binomial distribution with hyper-parameters $p$ and $q$. Thus, we have:

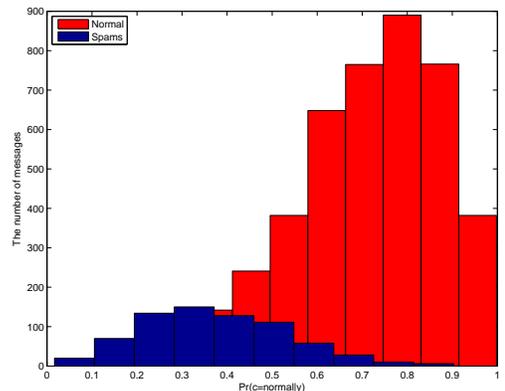$$c \sim p(c|\lambda) = Pr(c|1, p) = p^c(1 - p)^{1-c} = p^c q^{1-c}$$

After classifying the $i$th sample message, we generate the expected probability (this is the filtering output):

$$
\kappa \sim p(\kappa|c, \lambda) = \begin{cases} p(\kappa|c = normal, \lambda) &= \mathcal{B}(\kappa; \alpha_1, \beta_1) \\ p(\kappa|c = spam, \lambda) &= \mathcal{B}(\kappa; \alpha_0, \beta_0) \end{cases}
$$

Here, $\mathcal{B}$ represents beta distribution and its hyper-parameters are set as follows: $\alpha_0 = 3$, $\beta_0 = 5$, $\alpha_1 = 5$, $\beta_1 = 2$. Both thresholds ($h_1$ and $h_2$) vary between 0 and 1 by 1/30.



(a) N messages



(b) Distribution

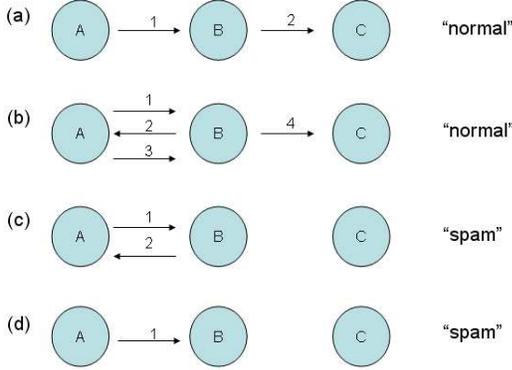**Figure 3: Displaying $\kappa = Pr(c = normal|y)$ for $N$ messages: spam (14.57%) and normal (85.32%)**

We have built an artificial dataset based on a Spanish database [7] which shows the proportion of spam as 14.57% and normal as 85.32%; that is, $q = 0.1457$ and $p = 0.8532$. The generated data have been plotted in Fig. 3-(a). A red cross represents normal message and a blue circle represents spam. This colouring scheme is also used in Fig. 3-(b). The graphs show that there are a lot of overlapping labels between 0.2 and 0.8. This overlapping section is considered as the uncertain region. Note that the challenge-response is not perfect and some of the spam might bypass the filter with correct responses, and normal messages might be filtered with incorrect ones. To model this imperfection, we use $e_1$ and $e_2$ to represent the ratios of False Positives (FP) and False Negatives (FN) in the uncertain region.

### 4.2 Traffic Usage Comparison

We have simulated and analyzed the traffic usage using the variable thresholds. Our framework considers three major stakeholders (see Fig. 1): a message sender (A), a message center (B), and a receiver (C).

First, we calculate the traffic used by an existing filtering method. Only the messages with filtering probabilities higher than the threshold $h$ reach C via B; other messages are deleted at B (only A to B). Suppose that $\mathbf{y}_{\bar{c}=type}^h$ for $type \in \{normal, spam\}$ denotes all messages filtered as $type$ in terms of $h$, the total amount of traffic used is the sum of $|\mathbf{y}_{\bar{c}=normal}^h| \times 2$, and $|\mathbf{y}_{\bar{c}=spam}^h| \times 1$ where $|\cdot|$ represents

the cardinality of a set: $N_{FilteringOnly} = |\mathbf{y}^h_{\tilde{c}=normal}| \times 2 + |\mathbf{y}^h_{\tilde{c}=spam}| \times 1$. This is because a normal message has two paths (A→B and B→C) but spam only has one path (A→B). In contrast, our hybrid model divides the mea-



**Figure 4: Four possible pathways for the hybrid method**

surable space into three different areas using two thresholds: $h_1$ and $h_2$. As a result, we have two more parameters to estimate: the traffic used by normal messages ($N_{un}$) and spam ($N_{us}$) in the uncertain region. Let $\mathbf{y}_{\tilde{c}=type}$ for $type \in \{normal, spam\}$ be a set of messages that have label $type$ as a ground truth.
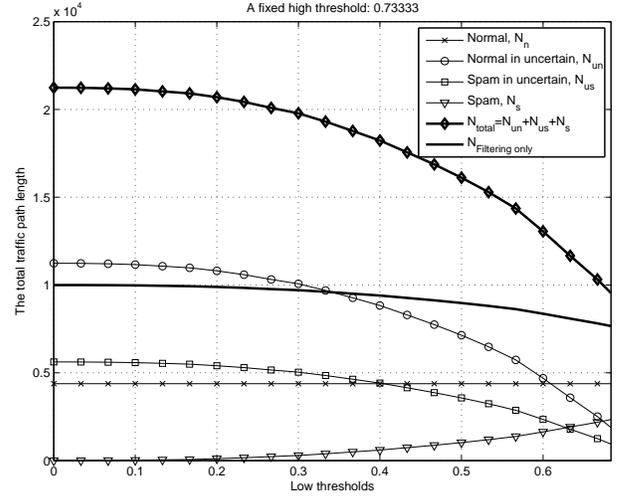
As Fig. 4 shows, there are four possible pathways. Firstly, a message classified as normal using the higher threshold is sent directly to C via B; the number of paths taken is two ($A \to B \to C$). Secondly, a message is in between the higher and the lower thresholds; a correct response is submitted by the sender and the message is classified as normal; the number of paths taken is four ($A \to B \to A \to B \to C$). Thirdly, a message is again in between two thresholds; this time no response is returned and the message is classified as spam; the number of paths taken is two ($A \to B \to A$). Lastly, a message classified as spam using the lower threshold is deleted at B; the number of paths taken is one ($A \to B$).
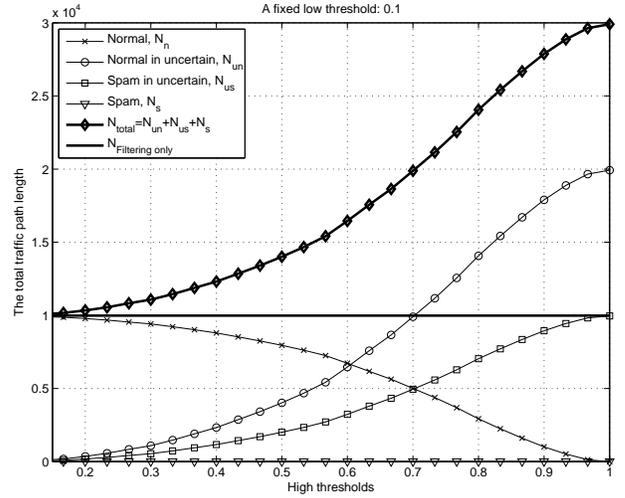
The traffic usage is calculated using:

$$
\begin{aligned}
N_n &= |\mathbf{y}^{h_2}_{\tilde{c}=normal}| \times 2 \\
N_{un} &= |\mathbf{y}^{h_1}_{\tilde{c}=normal} \cap \mathbf{y}^{h_2}_{\tilde{c}=spam} \cap \mathbf{y}_{\tilde{c}=normal}| \times (1 - e_1) \times 4 \\
&\quad + |\mathbf{y}^{h_1}_{\tilde{c}=normal} \cap \mathbf{y}^{h_2}_{\tilde{c}=spam} \cap \mathbf{y}_{\tilde{c}=spam}| \times e_2 \times 4 \\
N_{us} &= |\mathbf{y}^{h_1}_{\tilde{c}=normal} \cap \mathbf{y}^{h_2}_{\tilde{c}=spam} \cap \mathbf{y}_{\tilde{c}=spam}| \times (1 - e_2) \times 2 \\
&\quad + |\mathbf{y}^{h_1}_{\tilde{c}=normal} \cap \mathbf{y}^{h_2}_{\tilde{c}=spam} \cap \mathbf{y}_{\tilde{c}=normal}| \times e_1 \times 2 \\
N_s &= |\mathbf{y}^{h_1}_{\tilde{c}=spam}| \times 1 \\
N_{hybrid} &= N_n + N_{un} + N_{us} + N_s \quad (3)
\end{aligned}
$$

where $e_1$ and $e_2$ are the probability which normal people cannot respond. Here, $e_1$ and $e_2$ are fixed to values 0.02 and 0.01 respectively.

Figures 5-(a) and 5-(b) show the inner sections of the graph. The higher threshold is fixed to 0.73333, only the lower threshold increases from 0 until it reaches this value. The traffic is used less when the lower threshold increases. We have also monitored the traffic usage with the lower threshold fixed to = 0.1, and with the higher threshold in-



(a) A high threshold



(b) A low threshold

**Figure 5: Slides of an axis (with fixed threshold)**

creased from this value to 1 (see Fig.5-a). The traffic usage does not change with the filtering approach because the lower threshold is the same as $h$. As the number of messages in the uncertain region increases so does the traffic usage.

## 4.3  Variant proportion of spam

We have fixed the proportion of spam to 14.57% and normal messages to 85.32%. In this section we show how the performance is affected when these proportions change.

Table 1 describes a small number of samples from the nine different proportions. Each record has six different columns: proportion of spam (%), lower threshold ($h_1$), higher threshold ($h_2$), traffic usage (TA) of $N_{hybrid}$, ratio $\left(= \frac{N_{hybrid}}{N_{filteringonly}}\right)$, and accuracy $\left(ACC = \frac{TP+TN}{P+N}\right)$. It uses three different measures for the performance. If the traffic usage is less, we say the system is lighter and is more economical. The ratio is only close to 1 if the traffic used in the hybrid method is close to the amount used in the other. The accuracy measures the correctness of message classification. We can select practical threshold values for each

**Table 1: Traffic amounts and accuracy of hybrid methods in terms of thresholds**

| Proportion of spam | $h_1$ | $h_2$ | TA | Ratio | ACC |
|---|---|---|---|---|---|
| 10% | 0.1 | 0.2 | 10218.8 | 1.0239 | 0.91847 |
|  | 0.1 | 0.9 | 27698.96 | 2.7754 | 0.98312 |
|  | 0.4 | 0.6 | 13563.78 | 1.4218 | 0.95266 |
|  | 0.8 | 0.9 | 10493.06 | 1.6081 | 0.39682 |
| 20% | 0.1 | 0.2 | 10450.56 | 1.0487 | 0.83314 |
|  | 0.1 | 0.9 | 27859.08 | 2.7957 | 0.98391 |
|  | 0.4 | 0.6 | 13561.34 | 1.4659 | 0.94151 |
|  | 0.8 | 0.9 | 10050.6 | 1.5731 | 0.46933 |
| 30% | 0.1 | 0.2 | 10662.46 | 1.0707 | 0.74703 |
|  | 0.1 | 0.9 | 28114.76 | 2.8233 | 0.98462 |
|  | 0.4 | 0.6 | 13682.54 | 1.5179 | 0.94212 |
|  | 0.8 | 0.9 | 9404.9 | 1.5167 | 0.53119 |
| 40% | 0.1 | 0.2 | 10915.14 | 1.0972 | 0.6635 |
|  | 0.1 | 0.9 | 28188.34 | 2.8336 | 0.9853 |
|  | 0.4 | 0.6 | 13717.14 | 1.5641 | 0.93292 |
|  | 0.8 | 0.9 | 8721.52 | 1.4442 | 0.59632 |
| 50% | 0.1 | 0.2 | 11139.94 | 1.1214 | 0.58259 |
|  | 0.1 | 0.9 | 28341.5 | 2.853 | 0.98627 |
|  | 0.4 | 0.6 | 13523.66 | 1.5946 | 0.92748 |
|  | 0.8 | 0.9 | 8182.46 | 1.3902 | 0.66251 |

spam proportion to compare the performance. For instance, threshold values $h_1 = 0.1$ and $h_2 = 0.2$ can be selected in 10% spam proportion to show a reasonable performance of the hybrid method. However, if the system is concerned with achieving high accuracy and not with reducing the traffic usage, $h_1 = 0.1$ and $h_2 = 0.9$ values can be used. In a spam-dominant environment (for spam proportion of 50%), reasonable threshold values would be $h_1 = 0.4$ and $h_2 = 0.6$. Returning back to the figures for a spam proportion of 10%, $h_1 = 0.1$ and $h_2 = 0.9$ will be selected when accuracy is the most significant factor.

## 5. CONCLUSION AND FUTURE WORK

We proposed a hybrid spam filtering framework for mobile communication using a combination of *content-based filtering* and *challenge-response*. A message that falls into the uncertain region (after filtering), is further classified by sending a challenge (e.g. an image CAPTCHA) to the sender: a legitimate sender is likely to answer it correctly, whereas an automated spam program is not. Challenge-response protocols have been designed with the necessary cryptographic features. We have also shown the trade-offs between the *accuracy* and the *traffic usage* in using our framework. The simulation results suggest that, for a different level of spam proportion, the practical thresholds should be carefully selected according to the required level of the two measures.

In this paper, a synthetic dataset, as oppose to a real dataset has been used due to the following three reasons: firstly, we wanted to develop a generalized framework that is flexible and applicable to a wide range of applications (e.g VoIP spam filters[12]); secondly, it was not easy to find a real dataset since the challenge-response protocol is not a commonly used filtering method; lastly, this protocol involves a great level of human interaction and developing such a prototype (in order to generate our own dataset) was outside the scope. Our next step will be to contact mobile operators and forums like OMTP to collect real data, and evaluate our framework against other datasets.

Having the network operators charge for sending of text messages has been one of the big inhibitors to the growth of spam: even a cent per message might hugely alter the economics of a spammer. Assuming that a reasonable filtering method is in place, another hybrid potential is to force spammers to opt into a charging scheme where the cost of responding to a challenge is larger than sending an initial spam. For example, imagine it costs two cents to send a spam, then it would cost extra five cents to answer an image CAPTCHA.

## 6. REFERENCES

[1] P. He, Y. Sun, W. Zheng, and X. Wen. Filtering short message spam of group sending using captcha. In *Workshop on Knowledge Discovery and Data Mining*, pages 558–561, 2008.

[2] M. Healy, S. Delany, and A. Zamolotskikh. An assessment of case-based reasoning for short text message classification. In *In Proceedings of 16th Irish Conference on Artificial Intelligence and Cognitive Science*, pages 257–266, 2005.

[3] G. V. Cormack, J. M. G. Hidalgo, and E. P. Sánz. Spam filtering for short messages. In *In Proceedings of the 16th ACM conference on Conference on information and knowledge management*, pages 313–320, 2007.

[4] C. Dwork, A. Goldberg, and M. Naor. On memory-bound functions for fighting spam. In *In Proceedings of the 23rd Annual International Cryptology Conference (CRYPTO 2003)*, August 2003.

[5] R. J. Hall. How to avoid unwanted email. *Communications of the ACM*, March 1998.

[6] J. Golbeck and J. Hendler. Reputation network analysis for email filtering. In *In Proceedings of the Conference on Email and Anti-Spam (CEAS)*, 2004.

[7] J. M. G. Hidalgo, G. C. Bringas, Sánz E. P, and F. C. García. Content based sms spam filtering. In *Proceedings of the 2006 ACM Symposium on Document Engineering*, pages 10–13, Amsterdam, The Netherlands, October 2006. ACM Press.

[8] S. Shirali-Shahreza and A. Movaghar. A new anti-spam protocol using captcha. In *In Proceedings of the 2007 IEEE International Conference on Networking, Sensing and Control*, pages 234–238, 2007.

[9] L. von Ahn, B. Maurer, C. Mcmillen, D. Abraham, and M. Blum. recaptcha: Human-based character recognition via web security measures. *Science*, August 2008.

[10] W. Enck, P. Traynor, P. McDaniel, and T. Porta. Exploiting open functionality in sms-capable cellular networks. In *CCS*, Nov. 2005.

[11] D. Rogers. Mobile handset security: Securing open devices and enabling trust. OMTP Limited White Paper, 2007.

[12] N. J. Croft and M. S. Olivier. A model for spam prevention in ip telephony networks using anonymous verifying authorities. In *ISSA, New Knowledge Today Conference*, 2005.