

## Anti-Spam Techniques beyond Bayesian Filters

The most common spam counter-measurement today is keyword filtering. This method has many drawbacks. It is resource-intensive on the receiver side and often causes false positives or negatives. The email infrastructure as originally designed lacks a critical element: sender authentication. A receiving mail transfer agent (MTA) cannot easily find out whether a message is forged or authentic. This opened the way for spoofing and phishing.

### The Greylisting Method

*Evan Harris, 2003, <http://projects.puremagic.com/greylisting/>*

Greylisting is a general method of blocking spam based on the behavior of the sending server, rather than the content of the messages. Greylisting relies on the fact that most spam sources do not behave in the same way as «normal» mail systems. It looks at a triplet of information: the IP address of the sending host and the envelope sender (MAIL FROM) and recipient (RCPT TO) addresses. It rejects sending attempts from unknown «triplets» by replying with a temporary failure code.

The possibility of temporary failures is built into the core spec of SMTP (RFC 821). Any well-behaved message transfer agent will attempt retries if given a temporary failure code for a delivery attempt. This blocking comes with a minimal price of local resources. There is no unnecessary network traffic caused by Greylisting other than that associated with the connection itself. Since the receiving MTA does not check the content of the message at all there is very little processing overhead.

Greylisting is still quite effective against spam and virus mails. But spammers will surely adapt their software within short time when more providers implement Greylisting. The biggest disadvantage of Greylisting comes from the fact that users can't be sure that their email is received within minutes. Sometimes it can take hours until a message is delivered. For example if a sending server cluster attempts to send the same message from different IP addresses. A possible solution could be to look at SPF records that list the sending mail servers for a domain.

### Sender ID Framework

Big players in the email industry including Microsoft and AOL collaborate to specify a framework that allows receiving MTAs to authenticate senders of messages. This framework is thought to form a basis for complementary technologies like reputation and accreditation systems. Together, these technologies make it possible to build a spam-free layer on top of the existing email system.

Sender ID is a merger of proposals including SPF (influenced by RMX and DMP) and PRA (part of Caller ID). Service providers publish the IP addresses of all MTAs that are allowed to send email for their domains in DNS. Receiving MTAs check those entries. They extract domain names from the MAIL FROM line or by following the rules of PRA.

### RMX (Reverse MX)

*Hadmut Danisch, 2003, <http://www.danisch.de/work/security/antispam.html>*

RMX defines DNS entries that list all sending MTAs for a domain. A receiving MTA places a DNS query giving a domain name and gets back the list of IP addresses or hostnames that are allowed to send mail. It checks if the IP address of the sending MTA is in that list. Problems: RMX entries are static and potentially large.

### **DMP (Designated Mailer Protocol)**

*Gordon Fecyk, 2003, <http://www.pan-am.ca/dmp/>*

DMP adds for each IP address that is allowed to send email for a domain a separate DNS entry. A receiving MTA places a query with the IP address of the sending MTA and gets back a positive or negative result. Problems: it's impossible to define IP ranges and DNS caching is inefficient.

### **SPF (Sender Policy Framework)**

*Meng Weng Wong, 2003, <http://spf.pobox.com/>*

SPF combines the features of RMX and DMP by defining rather complex syntax rules. Domain administrators can choose which approach works best for them or they can even use both at the same time and even add rewrite rule entries for full flexibility.

SPF uses the MAIL FROM line to find out which domain is the origin of a mail message. This works well for mail that is sent directly and not forwarded. When mail is forwarded by a relay system problems with the return path arise. They can only be eliminated by redefining the MAIL FROM line. The great thing about SPF is that it can authenticate a sender before the mail message is received. Network traffic can be reduced this way.

### **PRA (Purported Responsible Address)**

*Microsoft, 2004, <http://www.microsoft.com/mscorp/twc/privacy/spam/senderid/default.aspx>*

PRA parses message headers and tries to find out the entity most recently responsible for injecting a message into the email system. From, Sender, and Resent-From headers are taken into account. The domain name found by PRA can be used to query a SPF entry. PRA works well also with forwarding relay hosts and mailing lists as long as they add appropriate header fields.

Disadvantages of PRA: the whole message has to be received before mail headers can be checked, and, most important, PRA is patent-protected by Microsoft. Developers that want to implement PRA in their software have to get a license from Microsoft. The license is provided free of charge. But open source organizations have repeatedly pointed out that these terms are not acceptable because they are not compatible with common open source licenses as those from Apache or GNU. The Apache Group for example would only agree with a standard that uses classic SPF but doesn't include PRA.

### **Domain Keys**

*Yahoo, 2004, <http://antispam.yahoo.com/domainkeys>*

Domain Keys signs the body and some message headers. It puts the signature in a DomainKeys-Signature header. Keys can be self-signed, as in PGP, and published in DNS following a decentralized, opportunistic encryption model. If a message fails signature verification, it should be rejected by the receiving MTA during an SMTP session, but in practice will probably result in some sort of warning sign in the MUA.

### **Reputation Lists**

Authentication of mail senders allows to prevent spoofing and phishing because users can trust in email addresses that are displayed in their mail user agents (MUA). But authentication itself does not protect users against spam because also spammers can add SPF entries for their domains.

Reputation systems are a key component of the future email infrastructure. They will help receivers decide if a mail from an authenticated sender is desirable or undesirable. Unfortunately I couldn't find concrete examples how these reputation systems should work. There are still only vague ideas. Reputation lists will enumerate domain names that are considered good. But which institution will set them up and step into new business segments? As with every presented proposal for authentication, reputation lists make sense only if leading email providers contribute to them and implement them, too.

All proposed changes to the email system cannot take place from one day to another. Transition periods are needed because mail server software and DNS entries have to be updated. Influential leaders have to do a coordinated first step. It seems questionable if the Internet community will accept the requirement to use Sender ID with all its implications if no widely accepted standard exists.