

THESIS

ANTI-SPAM MASTERS PROJECT

**Centrum voor
Wiskunde en
Informatica**



**Universiteit van
Amsterdam**



**AUTHOR:
DATE:
STUDY:
PROFESSOR:
COUNSELOR:**

**ING. SHAHJIHAN A. CHAUDRY
26-08-2004
MASTERS SOFTWARE ENGINEERING
PROF.DR. PAUL KLINT
PROF.DR. SJOUCHE MAUW**

Disclaimer to any solutions proposed; "The perversity of the Universe tends towards a maximum"- Finagle's law of Dynamic Negatives.

Acknowledgements

First and foremost, all praise to God almighty.

Most of the work that went into the thesis was not done by me, when you think about it I just structure everything and put my opinion to the matter. Hopefully people will still think I am a smart guy regardless of this point. Therefore I would like to thank the following people that helped make this possible:

My parents for supporting me during all these years, making me the man I am today and giving me the chance to study. My brother Usman and sister Jasmien, for caring.

My wife Naila, for her support and patience with my constantly changing moods, lack of attention for her, my obsessive need to study and a lot more

Sjouke Mauw, Paul Klint, Mark van den Brand and Alban Ponse for giving me the possibility to work on a project that seemed almost too much and being so very understanding and enthusiastic when times were hard.

Hans, Ravin, Sandra, Idris, Wilfred for reading my paper or offering to do so and lending an ear to my worries.

Everyone I didn't mention here, you know what you did and you probably know I appreciate you and what you have done.

ABSTRACT

Unsolicited mail (spam) is the number one threat to the use of e-mail. With the destruction of e-mail a part of the Internet will die with it. Spam has increased tremendously in the last few years. If this rate of growth will stay steady spam will outnumber solicited mail (ham) 9:1 within the next few years [1, 57].

The cost of spam can be measured in lost human time, lost server time and loss of valuable mail. Further more there is a level of annoyance at receiving a lot of spam. The dangers of spam are also on the increase, with the increase of virus-infected mail.

A general definition of spam does not exist because spam is different for every user. The definition of spam for the thesis can be seen as: Unsolicited commercial e-mail [2, 11].

The Objective: Create a solution that will run parallel to the current e-mail service and does not try to influence standardization.

By doing this, the solution should be able to stop 95% of all spam. 95% is the threshold amount of spam that the spammers cannot amount to lose [2]. With just 5% of their spam arriving as an advertisement, the cost of sending spam will be to high.

The Research Question: How could an enhancement to the current e-mail service be made, that will make it spam free and requires only standard technology and a little bit more user involvement of the communication partners?

The literature survey is the core of this project and has led to several solutions; these are analyzed and theoretical solutions are discussed. Currently it is possible to stop 99.984% of using the available spam existing solutions [3]. Filters have evolved into the primary means to stop spam. They can be complemented by several solutions like black/white-lists.

Until recently there were no organized efforts to stop spam. A few organizations have now been founded similar to the Anti-spam research group (ASRG). These organizations provide a framework that individuals can use to take their innovations to the next level.

Principles that help the further understanding and development of spam solutions are: changing spam [19], false positives [18], nature of digital technology [20], enough technology [3], eye-space [30] and shared ownership [6].

The diversity principle is a theory that explains why it might be possible to stop spam. It states that each person has a different set of preferences and should use differing spam solutions, which suit their mail usage. By combining these preferences over a large group of people a broad scope of spam definitions are created and a diverse approach to stopping spam [4]. Therefore when a spammer is targeting a group of people and there is a lot of diversity he will not be able to find a single method or single message that will fool everyone.

Techniques that are still theory and have to be tested like inoculation, just in time filtering, dynamic mine fielding and tier foldering give hope for a level of accuracy in the vicinity of 99.9999%. Taking spam solutions to the next level.

The future holds 3 prospective realities:

- 1) Worst case – 9 out of 10 mail is spam, solutions cost money and don't work properly, user privacy is lost, false positives are common. >> The death of e-mail.
- 2) Most probable – E-mail stays free, solutions will not drain funds, 1 false positive per user lifetime, privacy of users is maintained, spam will practically be stopped. >> E-mail stays useful.
- 3) Ideal – E-mail stays free, solutions don't cost anything, no false positives, user

privacy is maintained. Because the spammers, saw the error of their ways. They all decided to join Greenpeace and work for a better tomorrow. >> E-mail is better than ever.

Conclusion

Spam is a problem that costs a lot of money and needs a solution. During the project a lot of information about spam and a lot of solutions for spam have been found. Even though there are several good solutions available spam still populates the inbox of many mail users. This is because the solutions rely on the user be very actively involved with the solutions. The filters and black/white-lists other solutions, although quite good, are not strong enough to attain a 99.9999% accuracy level. A lot of users do not want to use the existing spam solutions because they do not find them accurate enough, generating too many false positives. Most users that do use the solutions do not attain the top ratings of 99% accuracy that the developers themselves do project. The levels the users do attain are apparently not enough to halt the current spam threat.

To possibly stop spam solutions are needed that can reach the 99.9999% accuracy level when pressed and can reach a 99% level without much trouble. The second condition is to accommodate the normal mail user that cannot invest too much effort and time.

To create such a solution a collection of existing and theoretical solutions will have to be combined. The synergy that results can theoretically accomplish the requirements.

The experts are focusing on developing techniques that can detect the spam at the receiving end of the transmission. Other solution categories usually are theoretically valid but have been deemed ineffective in real life tests.

The ultimate objective of an absolute ban of spam is not a real prospective. An analogy; we can secure a bank as much as we want, but people are still going to rob it. "*Because that is where the money is*". The people that will probably keep spamming are the ones that have very low cost and can live off a very low profit margin.

CONTENT

ABSTRACT	I
CHAPTER 1 BACKGROUND	2
BACKGROUND INFORMATION	2
CONTEXT OF RESEARCH	2
RESEARCH QUESTION AND OBJECTIVE	2
CHAPTER 2 PROJECT PLAN	3
ACTIVITIES	3
DELIVERABLES	4
CHAPTER 3 EXECUTION OF THE PROJECT	5
PROCESS	5
LITERATURE SURVEY	5
CHAPTER 4 ANALYSIS	7
PROBLEM ANALYSIS	7
EXISTING SOLUTIONS	9
POSSIBLE SOLUTIONS	12
THEORETICAL SOLUTIONS	13
ALTERNATIVE SOLUTIONS	14
CHAPTER 5 CONCLUSION	16
CHAPTER 6 EVALUATION	17
POSITIVE	17
NEGATIVE	17
REFLECTION	17
SELF-ASSESSMENT:	17
REFERENCES	18
FIGURE LIST	20
APENDIX-A LITERATURE SURVEY FINDINGS	21

BACKGROUND INFORMATION

At the moment of writing, the major part of all e-mails that are received are unsolicited (spam). There are now various other forms of spam; Usenet spam, SMS spam, web log spam, etc.

Spam is the primary threat to the survival of e-mail as a useful communication medium. Figure 1 gives an overview of the amount of spam a user received between the years 1997 and 2000 [55]. The attack of spammers has increased dramatically during the past few years and has consequently decreased the usefulness of e-mail. We are not only experiencing annoyance but are losing millions in the form of human resources and server time. Fighting spam is appropriately called a war. In our efforts we are spending (or rather wasting) time and money trying to stop spam, which could be employed for any number of constructive projects.

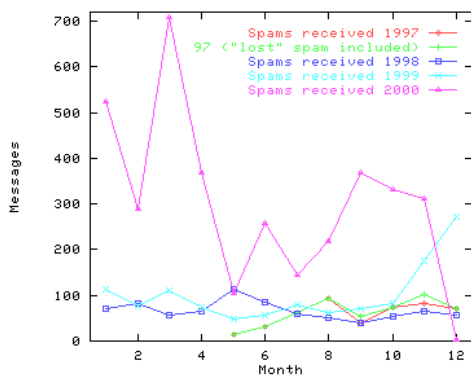


Figure 1. Spam in Europe 1997-2000

CONTEXT OF RESEARCH

At this time several organizations are being formed to fight the war against spam. Organizations make it possible to create solutions in an structured manner. A large amount of software is being developed to stop spam, of which the majority are filters.

Despite the availability of solutions to spam, users are reluctant or unable to use them. This is primarily due to the lack of transparency and relatively difficult use of the solutions. The training of for example

filters is beyond the ability of many ordinary users.

Some of the ideas that are circulating at the moment do not seem to be appropriate for implementation. For example a micro payment system that makes the sender pay for sending a mail. If the receiver accepts the mail as being solicited, the money will be refunded. If it happens to be spam the spammer has to pay the bill. While this sounds ideal, the reality is that there is currently no micro payment system that works. There definitely isn't one that could be used on a global scale.

The relatively new approach to the issue is persecution of the spammers and the organizations that fund their activities by using the law. The persecution of spammers by the law is not very effective. This is mainly because the spammers are crooks. They will try to avoid being caught. When they are caught the law is not clear enough for them to be prosecuted. For more information on the legislation and developments in this area, see <http://www.isipp.com/>

RESEARCH QUESTION AND OBJECTIVE

The Research Question: How could an enhancement to the current e-mail service be made, that will make it spam free and requires only standard technology and a little bit more user involvement of the communication partners?

The Objective: Create a solution that will run parallel to the current e-mail service and does not try to influence standardization.

CHAPTER 2 PROJECT PLAN

The project plan is coarse grained due to the nature of academic research, which has a high level of uncertainty. The planning leaves room for change and redirection.

ACTIVITIES

During the project there was a subdivision of core activities and supporting activities. Outlines of the activities and their relative worth to the project are stated. The worth of the activity is stated between the [...] brackets. Scale: Very Low – Low – Medium – High – Very High

1) **Planning the project**, completed in the first week of the project and updated with information of activities that were (not) performed. [MEDIUM]

Contents:

- o Determine scope of project;
- o Specification of sub activities;
- o Planning of effort and time spenditure;
- o Specification of products;
- o Risk analysis and mediation;
- o Methods needed per activity and product.

2) **Literature survey**, the information needed to successfully execute the project was collected during the survey. This activity was executed parallel to almost all others during the project. [HIGH]

Contents:

- o Determine scope of survey according to scope of project;
- o Determine hot spots;
- o Collection of documents;
- o Examination;
- o Asses usable literature.

3) **Analysis of problem and domain**, with the aid of the literature gathered up to this point, the problem was analyzed in depth. The domain knowledge that had been attained in the literature study was tested against the ability to independently describe the main characteristics. This was then checked, to be in accordance with the established authorities. A start was made

towards documenting the solution, mainly in the form of sketches and notes. [HIGH]

Contents:

- o Brainstorm;
 - o Determine scope of think tank and domain;
- o Management Analysis;
 - o Analyze and select possible solutions;
- o Determine relevancy of literature to solution;
 - o If literature is not up to speed, assert extra effort to survey new literature;
 - o Document the first draft of the solution.

During this phase testing domain knowledge showed certain defects in the area of mathematic references and conventional e-mail transmission. This led to a further extension of the survey in this area.

4) **Solution design**, the spam solution that was devised came in the form of suggestions. The conclusion was made that the current technology is sufficient to provide a proper defense against spam. The need to develop new techniques was reduced. Instead the form and method of using existing techniques was now the main focus. Various known solutions were scrutinized and in some cases successfully falsified. [VERY HIGH]

Contents:

- o Analyze the core functionally;
- o Asses solution strengths and weaknesses;
- o Falsify if possible;
- o Propose theoretical solutions.

A detailed plan was not forthcoming because the solutions are very conceptual; it is hard to decompose them to low levels of detail.

5) **Writing thesis**, the thesis will be written parallel to all the other activities during the project. [VERY HIGH]

Contents:

Main

- o Determine goal of thesis and way to achieve the best relay of project message;
- o Writing the thesis;
- o Write main body of thesis;

Analysis – Creation of comprehension

Design

Thesis

Management –
Project start-up

Analysis – Information
gathering

- o Write abstract;
 - o Write introduction, thanks and conclusion;
- Details*
- o Make inventory of quotations, definitions, diagrams, illustrations;
 - o Write bibliography / list of references relevant to thesis;
- Check*
- o Check for unintentional plagiarism;
 - o Determine the structure and typographic design;
- Final*
- o Review of thesis by several third parties;
 - o Create enough copies.

DELIVERABLES

Deliverables that will result from the activities are:

- 1) Research formulation;
- 2) Project plan;
- 3) Literature survey report;
- 4) Solution design report;
- 5) Mockup;
- 6) Thesis.

CHAPTER 3 EXECUTION OF THE PROJECT

PROCESS

The planning was a ruff estimation of the time and effort that would need to be spent to achieve the desired goals. This was done because there was not enough experience with projects of such a short time period and the uncertainty in academic research. To try and keep productivity high, time boxing was proposed. But this would not be a viable approach because the course the project would sail was not clear.

When trying to create a new way of handling spam a few preparatory steps need to be made. One must know what technologies already exist and if they are in anyway successful. All the ways that an idea is formed must be done in a scientific manner, so a justification can be given.

The following steps were used when trying to make a spam solution.

1. Problem definition;
2. Problem scope;
3. Solution definition;
4. Solution scope;
5. Document the principles and the proper way to use the solution;
6. Try to falsify the solution.

LITERATURE SURVEY

After examining the difficulties in the project it became obvious that it would lean heavily on the literature survey. The objective of the literature survey was to ascertain the state of current affairs in the field of spam development and the attainment of knowledge about spam.

The literature survey claimed a lot more time than was anticipated. This was not necessarily a bad thing. Because one of the conclusions that could be made after the literature survey was that a lot of work has already been done to fight of spam. This conclusion and the subsequent analysis of the spam technology and methods prevented the reinvention of existing solutions.

SOURCES

The literature survey was mainly be conducted on the World Wide Web. The www resource provided a sufficient amount of literature on the field of interest. Next to the www the libraries of some educational facilities were consulted for previous student research and publications on the subject. All most all the literature was found on the Internet. While searching the Internet only one definite reference was made to a book. This was an indication that the main body of knowledge about Spam resides on the Internet. Other relevant institutions are companies that engage in this field of expertise and governmental agencies.

To explore the web a utilities will be used in addition to browsers. The utilities are the Kazaa utility, e-donkey network and the e-mule utility. A small search for other useful utilities will be done. With these capabilities it was possible to find enough literature for the objective of the project.

The four perspectives that will be used are:

1. Spammer: sources that describe the way Spammers go about their work, written by the Spammer himself.
2. Spam solution creator: the various solutions that are available at the moment, written by all parties.
3. Current mail service: a description of the way the e-mail works, written by all parties.
4. Future development: possible solutions to the Spam dilemma that are being developed, written by all parties.

CONCLUSION

A lot of information and solutions for spam have been found. Yet there spam still seems to get into the inbox of many mail users. This is partly because the user does not employ every method to fight spam; usually a user only has a basic spam filter or a block list. This is not enough to minimize the current threat of spam.

The objective to achieve an absolute ban of spam is not a real prospective; a minimization of the amount is the best that can be achieved. To accomplish this the spammer needs to be waylaid on several levels. Level indicates the distance from the user, where high is far from the user and low is close to the user.

1. Server-side, this is a high level. The goal on this level is to minimize the amount of spam that can be sent, so the user has less to filter through.
2. Client-side, this is from middle to low. The goal at middle level is to communicate with peers about know spam and eliminating these, thus reducing the amount the filter has to process. The low level goal is the filter, which intelligently quarantines the amount that is still suspected spam.

This leads to the situation that the user has a small quarantine box, which is easy to check and less likely to let false negatives go by.

Various organizations are competing for the position of spam leader, with software or standards that will enforce this. However the experts say that the modification of the standard is not the way, because this usually means the killing of the free and easy communication via e-mail. They instead point in the direction of an eclectic solution that will build upon the strengths of each method.

PROBLEM ANALYSIS

What is spam: Spam for the purpose of this thesis is defined as unsolicited commercial mail [7]. It is necessary to point out that some people want to receive these messages. There is an audience for e-mail advertising, regardless of the product that is being sold. Spammers are trying to reach these people. Spammers do not know which people comprise this group. To reach them, they send spam to as many people as possible. This is done because they don't know who will respond to the message and who will not.

Spammers are generally technically skilled individuals that are hired by companies to send spam. By using a third party, the companies try to keep themselves from getting sued [8]. For a company spamming can be very lucrative if done right. For example a company is selling wonky dolls for 50 dollars a doll. If the company lets the spammer send out 10 million mails and the response rate is just 0.1% it will make half a million dollars [9].

Spammers get e-mail addresses by foraging them from websites, newsgroups etc [10]. It is possible to turn this into an advantage, by fooling spammers with fake e-mail addresses and thus harvesting their spam. This will be discussed in the section theoretical solutions.

A definition of spam in the free online encyclopedia Wiki [11] states that:

“Spam typically refers to any of the following:

- *Spam - a brand of canned meat sold by Hormel*
- *Spam (Monty Python), a comedy sketch involving this meat, which in turn gave rise to the phrase:*

- *Spamming - unsolicited electronic messages, such as spam (e-mail);*
“ [11]

Spam is not limited to e-mail. Spam exists in text messaging services (SMS), newsgroups and other communication media. In the case of SMS, spam can cost even more than it would when received through mail. For example, a user has subscribed to receive a notification via SMS when she receives e-mail at her mail account. She will have to pay for every SMS received regardless if it announces a spam or a ham [12].

That spam is different for everyone creates a problem when trying to train filters. Training the filter is something that has to be done by the user himself. For example users named Bob and Tom. What Bob considers to be spam is not the same thing Tom considers to be spam. So when Bob trains his filter, Tom cannot use it because it does not know what Tom considers to be spam.

What is spam conclusion: The definition of spam for the thesis can be seen as: Unsolicited commercial e-mail [2, 11]. This definition is by no means absolute.

Why is spam a problem?

Spam does cost money [13]. There are three types of cost: capital, staffing and business [14]. The users lose time and various ISP's lose money, trust, working hours or even operation of their servers due to spam [15, 16]. It is estimated that spam might cost the billions of dollars in the near future [17]. Research has not only given these indications, ISP's have indicated the amount of money they are losing [14]. The other reason why spam is costing money is because of the loss of valid mail, by losing it in the flood of spam. This can result in the loss of business, as mail is an important form of

communication for various businesses [18].

Changing spam: It is generally accepted that spam changes rapidly and is highly unpredictable but recent research has shown that spam actually is quite predictable. [19]

“Analysis of a high-fidelity statistical model of historical spam characteristics reveals long periods of minimal variation, occasionally “punctuated” by brief bursts of sudden change. “

While spam may stay constant for a longer period of time than expected, the changes it undergoes are still unforeseen. These changes are an innovation in hiding spam that misleads spam identification technology into thinking it is ham. At the time of change, the filters will not be able to stop spam from breaking through. Thus rendering them ineffective until an adaptation can be made, that can deal with the change.

Changing spam conclusion: Spam has time stable characteristics that do not change. These time stable characteristics enable the prolonging of filter life spans and other identification technology. The short sudden changes should be followed by a swift reaction of the anti-spam community. This indicates that whatsoever measure we take spammers might always find a weakness in it.

Enough technology: Developers in the field of spam solutions have reported their accomplishments in stopping spam [20, 3]. They indicate that stopping 99% (and more) of spam is possible. If it is known how to stop 99% of spam, why is it not being applied?

The main reason is that it is still hard for a user to use a spam solution. In most solutions some form of filter is used to identify a spam message. A filter has to be trained, so that it becomes efficient in recognizing spam and ham [21]. The user usually trains the filter [22]. When training the filter, it is supplied with information about the preferences over an extended period (one or more months). It seems a

user-unfriendly process, because it takes about a month, to train the filter before it becomes truly effective. On the assumption that a typical mail user is not a very technical person nor does he/she want to spend too much effort in making an application work. It might be stated that user-friendliness and transparency of technology that is used in the application is still an issue for filters.

Enough technology conclusion: The existing technology has enough capabilities to stop a great deal of spam, but user-friendliness needs to be improved.

False positives: The term, false positives, refers to a legitimate message that spam solutions wrongly see as spam. Possibly the most important problem in fighting spam is the occurrence of false positives. False positives are generally an unacceptable error for users. This is because of the loss of important mail messages. The cost of false positives for the US economy is rated at 3.5 billion US dollars in 2003 [18]. While the cost of spam is estimated at 8.9 billion US dollars, false positives can also be very harmful for business relations [23].

False positive conclusion: False positives are the worst side effect of spam solutions. It costs a lot of money and a very low level (1 per 2000) of false positives is already too much [24].

Nature of digital technology: The nature of digital technology lies at the heart of spam. Spam is a profitable business because the cost of reproducing spam is virtually zero [25, 9]. For most spammers Internet is partly seen as a place where one can anonymously surf the web and send mail.

There are several different ways to send spam, directly through a trial account, through an intermediary system, via an open gateway located in some distant country, etc [26]. These options make it very hard to identify and catch spammers.

Nature of digital technology conclusion: Due to the nature of digital technology it is

Weakness of spam

Strength of spam

Weakness of spam

Strength of spam

difficult to trace the spammer [27]. Therefore interception of spam is a more viable solution.

Shared ownership: Spammers and their supporters control a part of the Internet, i.e. they own gateways, servers etc [6]. One of the reasons spammers are so successful, is that they have the ability to continue their work. With the available open gateways and other facilities they can relatively easily send spam. For example the Internet users in china, which is the number one producer of spam [28], will probably not be as easily cleansed of open gateways as the U.S.A.

Shared ownership conclusion: Spammers will be able to send spam, because certain types of facilities (like open gateways) cannot be wholly abolished. This must be anticipated and frameworks that can reduce this to a minimum are CAUCE and governments [29].

Eye-space: Although spammers may do a lot to disguise their spam from filters, they cannot disguise their sales pitch without distorting it enough to render it ineffective [30]. Every other area of a message can be modified extensively; they do not have to adhere to the eye-space rule. Because the target is human, spam eye-space will always have to be readable and understandable.

Eye-space conclusion: The main weakness is that it will always have to be understandable for a human. This gives us the possibility to strip the mail to its core message and filter on the eye-space content.

Other points of weakness: The definition of hop count is; "The number of signal regenerating devices (such as repeaters, bridges, routers, and gateways) through which data must pass to reach their destination". For spam that is being relayed the hop count is usually higher [31]. This is a lesser indication because it will only apply to circular relayed or extensively routed spam.

Bulk mail is an acronym for spam. Spam is usually sent in large numbers. This is necessary to reach the users that will respond to the message. Therefore the arrival of a large amount, at an ISP, of identical messages is a good indication of spam. The amounts are usually much larger compared to a mailing list.

Base64 encoding of the header part of a message is a strong indication that the message is spam [32]. This is done to increase the difficulty with which the spammer can be traced [33].

Other nuisances

When a user sends a message that the filter at the recipients end erroneously considered to be spam, the sender does not know that his message is identified as spam. Assuming everything went fine, the false positive is not pursued. The sender can send several following messages that can have the same result causing serious harm.

EXISTING SOLUTIONS

The current solutions to spam employ a variety of techniques. Some of the tools available like Spamassassin are very successful at stopping spam [34]. There still remains room for improvement [54].

8

One of the principles that are important to effective construction of solutions is the diversity principle.

Diversity: Spammers employ a lot of tricks to throw identification off course. These sometimes mislead the filters, especially when a spammer made up a new trick [33]. Some filters will let it through while others will pick up on it. Successful spam is something that requires a lot of effort from the spammer. The way spammers attack is by analyzing a specific filter type and exploit a weakness in its rules [21]. When the general population of users takes advantage of a range of different filters, it will not be worthwhile

for spammers to target one type. Targeting several types of filters is very complex. For the same reasons Barry Warsaw also recommended avoiding monocultures in spam solutions [4].

Diversity conclusion: Diversity in spam identification techniques will enhance the overall stopping power of spam identification.

Filters

The use of filters is the main form of defense against spam. There are two groups of filters, Heuristic and Bayesian [30]:

“The filters are called “heuristic” because they determine only the probability that a message may be junk e-mail, based on rules created from empirical observation of thousands of junk e-mail messages” [35].

A Bayesian filter needs no human intervention to generate the feature recognizers. By breaking the incoming text into words, each word becomes a feature [39]. The Bayesian filter counts the features (divided into spam and non-spam) frequencies and from that ratio determines a local probability on a message [56].

The latest range of filters decodes e-mail to its eye-space message before filtering. This use of eye-space is what gives filters a distinct advantage.

CRM114’s 99.87% is the highest pure filter accuracy that could be found [3]. This should be enough, but the existence of a false positive is a major downside to any filter. A solution is getting 99.9999% (one per million e-mails) accuracy. This equates to an average user (50 mails a day) receiving about 0.675 per 45 years of use. This type of accuracy (99.9999%) is called a Gaussian tail and is not within the reach of filters. At this rate a user receiving 50 mails a day, would get one false positive every 75 years.

Theoretically this level should become a reality when using inoculation and minefields, in addition to the standard

filtering, black/white-lists etc. techniques. This will be discussed in the following segment “*Theoretical solutions*”.

To achieve a level of 99% accuracy with the CRM114 a month of training is needed at a high rate (+100 messages a day) of live e-mail. Most users don’t have a high mail rate. Therefore the time span will be considerably longer than a month. This is a long time for the average user that is used to plug and play style applications.

A solution is training the filter with the aid of spam and ham archives also named corpus [35]. By using the corpus a high rate of incoming e-mail can be emulated, shortening the training period. The corpus should contain 5600 (200x7x4) messages. Theoretically a higher amount of messages should result in a better filter. A corpus that is designed for training purposes should contain a weighted cross selection of the spam or ham population [36].

The use of a corpus artificially creates a high amount of incoming mail, thus enabling the user to train the filter more than with his usual amount of incoming mail. This approach does not lessen the amount of work; it only shortens the time in which the filter is trained.

Figure 2, shows the progress a filter usually makes if trained correctly and consistently by one user. It also shows that the 0 for amount of false positives is not reached.

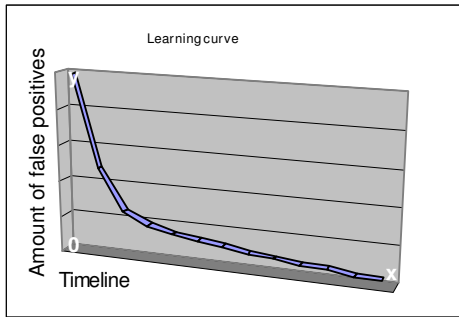


Figure 2. General filter error learning curve

Indications by ISP's have been given that they have archives containing several terabytes of e-mail [spam conference MIT]. Due to privacy reasons these cannot be shared, for the time being. The accumulation of ham is difficult because people are reluctant in sharing their personal messages. The collection of ham will have to be done by receiving donations.

The next step would be to develop an automated training system. Selection of pre-composed preferences could increase the transparency for the user. The underlying purpose is to lower the threshold for user involvement, resulting in less of the total amount of spam being read.

Black/white-lists

Black/white-listing is a common technique used to stop spam. Blacklists contain the addresses of spammers, when a message comes in. A check is performed to see if the sender is listed in the blacklist. If he is, the message is automatically treated labeled as spam. White lists are the opposite; they contain users that are verified contacts [25]. These users may send messages that seem spammy, but because they are listed on the whitelist, will be treated as ham.

Black/white-lists is a technique that, employed by a large system such as MAPS, has been shown to catch only 24% of spam, with a very high false positive

rate [22, 29] Reviews state that black/white-lists are often susceptible to misuse, because humans have to manage the lists of blocked IP-addresses. Furthermore a trend is that spammers are getting onto credited white lists, this is disconcerting because it means that they are able to fool humans.

The technique has come along since then and does help the decrease of false positives in filters like spamassasin and CRM114 [3, 34]. But it might be wise to only use it as a complimentary technique, as is done in the mentioned solutions.

Distributed hash databases

“Vipul's Razor, Pyzor and DCC are collaborative spam-tracking databases, which work by taking a signature of spam messages. Since spam typically operates by sending an identical message to hundreds of people, these databases short-circuit this by allowing the first person to receive a spam to add it to the database -- at which point everyone else will automatically block it” [34].

In the following, figure 2, the principle of distributed hash databases is shown for the Pyzor network [53].

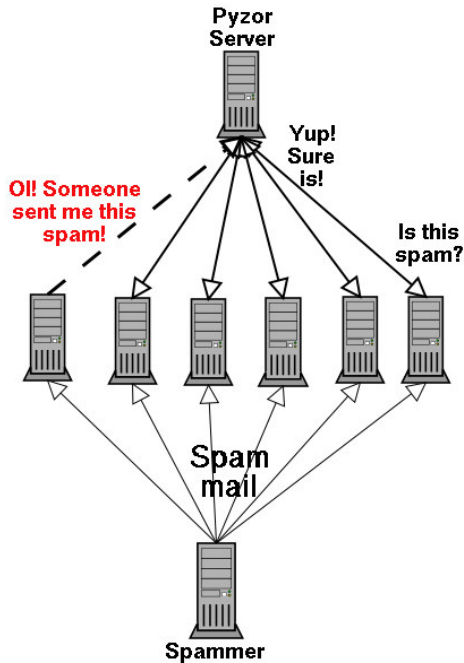


Figure 3 Pyzor schema [37]

This seems to be an ideal solution, with users informing each other. This is similar to inoculation that will be discussed in the section, theoretical solutions. The downside is that someone has to classify the e-mail to be spam; it is not automated and thus biased. Because of the human judgment element (similar to the problem in black/white-lists) the solution is open for corruption.

Furthermore it creates a digital signature of the body of an email message, instead of features or descriptive information. It then asks a central server if any one else has reported that digital signature as a spam mail [37]. If the message could be changed for every user, the system would no longer recognize it.

Due to the communication through a central server, the system will probably not scale well to a global level. The amount of cross-referencing required would increase in relation to the amount of users and spam. Such a degree of relational DB that communication with each other would require too much

resource and would probably be lacking in performance.

The arguments that are given are general arguments that can be applied to any system with central points of communication/DB and a large amount of users when scaling up to a global level

POSSIBLE SOLUTIONS

Sender filtering

Assumption, the spam identification tools the sender uses are similar to those the receiver uses.

The software can see a message as spam for many reasons, but the most likely are that a filter judged the message to be spam. When false positives occur the sender can take a preventative measure, by employing his own filter to indicate faulty net etiquette by using his own filter as a word checker.

Upon sending a message a check should be performed to see how the message scores. This is an indication to the sender that he should change his message like removing certain words that are high indicators for spam that will be highlighted by the application. The user that would benefit from this type of feedback is the typical user as mentioned earlier, that doesn't have a lot of technical knowledge or insight.

Feedback

Feedback is sent to the sender when mail is recognized as spam. This form of feedback can be used to attack a Bayesian filter as John Graham Cumming points out by using an "evil" Bayesian filter to repeatedly send spam with included random words [40]. Resulting in a distorted learning curve for the evil filter, which gets to know the key words that are needed for the good filter to classify a message as ham. What is remarkable is that it just needs one of the detected keyword to be added to the message for it to be classified as ham.

Graham-Cumming further more states that this form of attack can only be stopped when disabling all feedback measures. Leaving us once again with the lack of information for the authentic sender. Thus the authentication of the sender could be combined with other techniques for a better protection of the filter.

THEORETICAL SOLUTIONS

As was described in paragraph “*filter*”, 99.9999% accuracy is not within the capabilities of filters [30]. A combination of techniques must be used to achieve this level of accuracy.

Just in time filtering, Inoculation and Dynamic mine fielding are techniques that are still in development. The preliminary results are that they function well but the test data is not yet available for review [30]. Therefore only a theoretical idea of the techniques can be discussed.

De Sade observation – “*One man’s pain is another man’s pleasure.*”

With this in mind, one person who receives spam can let others profit from his pain.

Inoculation

Internal inoculation, when a group of users create a network where they inform each other about the spam they receive. The error rate of the system for discrimination of spam is related to the inverse of the number of users. With 10 users participating a 10-fold improvement in filtering is achieved [30].

However, the time it takes for one user to inform the others might be too long. The spammer could already have sent spam to all the users. To provide the needed time just in time filtering is employed.

Tier folders

The assumption behind tier folders is that when a false positive does occur, the mail

is weighed close to the border value of being spam [38, 39]

If this is the case, a folder system could be made where per percentage range a folder would exist. This would enable an easy retrieval of the mail. In contrast to the large spam folder that usually is created which makes it hard to find the false positive. The folder with the borderline percentage of mail would contain only a small number of messages. Thus making it easy to search through a spam folder of choice when searching for a false positive.

The possibility exists that a grievous error by the filter would go unnoticed because it would be placed in one of the lower level folders. The user would in a normal course of action not check the lower folders for false positives. Thus the false positives would go unnoticed.

Dynamic mine fielding

E-mail minefields form a defense against site-wide spam attacks; where all the mail accounts on an entire site are spammed in a very short period of time. The minefields are in essence large sets of bogus e-mail addresses. These are intentionally leaked to the spammer, and consist of improbable names. The spammers use tools to harvest the addresses, which do not discriminate against these names. Since no human would send a mail to the addresses, any mail arriving here is spam. Older versions of e-mail minefields are honey-pots and spam traps [41, 57].

When the message is received the users in the inoculation network can immediately be informed. When the information about spam is coming from the spammer himself, it is called External Inoculation [42].

More usefully, spammers usually attempt to falsify their headers and hide their IP addresses. However, during the SMTP transaction from the spammer to the minefield address, the spammer must reveal their actual IP address. The spammer cannot spoof this address, as the

SMTP transaction depends on at least the RCPT OK section of the transaction being delivered correctly to the spammer, and that can only happen if the spammer reveals a correct IP address during the socket setup phase.

At this point, the targeted site and any site cooperating with the targeted site can immediately blacklist the offending IP address. This blacklist can be either a “receive and discard”, or “refuse connection” situation [30].

Just in time filtering

To create the largest amount of time for inoculation to occur, just in time filtering can be used.

Current e-mail delivery systems filter their messages on arrival time (SMTP time). This does not allow enough time for inoculation to occur, in some cases. An extra moment of filtering, at the user-read-time (when the user is retrieving his mail from the server) provides adequate time.

Inoculation is partially being used in Dspam [43]. The implementation is not done in the same manner as described here, but it is an encouraging move towards implementing these solutions.

ALTERNATIVE SOLUTIONS

Pretty Good Privacy (PGP)

PGP is the most used encryption standard when sending secure mail [44]. With 128-bit encryption and digital signatures for secure email, the standard is still strong enough for most purposes. Encryption ensures that a message cannot be read by anyone other than the intended recipient. Digital signatures provide verification of the creator’s identity and that the message was not tampered with during transit [24].

Authentication is the most important part of PGP in relation to spam. Authentication is performed using certificates that are obtained at a Certification Authority (CA) like VeriSign and RSA. A flaw in PGP lies

in the ability of basically any organization to become a CA. Thus also granting certificates to spammers. This defect exists in the chain of trust that is not reliable [42]. So if you cannot trust the CA you cannot trust the certificate. Furthermore a compromised key that is spread over the Internet gives the spammers the ability to spam with a false identity. This form of misuse can only be remedied after the incident is reported, thus spam is not totally solved.

PGP works well for users on both the sending and receiving end, who are willing to use encryption when sending mail. The fact remains that not everyone, a user receives e-mail from will use PGP, but the user will still want to read this mail and will not automatically label it as spam if it is not encrypted. Therefore PGP cannot provide a solution to the spam problem, it is only appropriate when wanting to send secure e-mail.

Sender pays

Camram is a solution that is based on the sender pays concept. The idea is that the person sending the message has to pay in time or money. This will result in a small amount for the individual user but in a massive (too large to support) amount for a spammer [45].

Camram uses the term postage, which is a partial hash of the message. The postage is generated using Adam Back’s Hash cash [46]. Hash cash makes the user pay about 15 sec of CPU time by calculating n-bit partial hash collisions on the message.

Camram offers four levels of filtering:

- 1) Is the postage present?
- 2) Is the source email address present in a list of known addresses (white list)?
- 3) Is the message a response to a postage-due notice?
- 4) Does the message pass a spam discriminator?

The upside is a lower amount of false positives; the downside is that people can refuse to use postage. When a sender does not send postage (like a spammer) the

system will send a postage-due notice. However spam arrives like a tsunami, in a large amount and in a short interval at the ISP. When all recipients send postage-due notices this will result in a “*knock-on tsunami*” effect [41]. This effect could be crippling for the ISP, thus reducing the scalability and usability of the application.

The “*knock-on tsunami*” is also the downfall in systems like Spamarrest. They require a sender confirmation, that the sender is human [47].

Public Key Infrastructure (PKI)

The partial implementation of a PKI to stop spam will be discussed in the form of the Yahoo domain key system that is being in its finishing stages of implementation [48]. While this is very similar to PKI, it lacks a CA. This could be an essential flaw, because it makes it possible to receive a key while not having to identify oneself in real life to a CA. By keeping all the interaction on the web, it creates a weakness. By having to identify the user in real life the identity of the user is

positively ascertained with official documents, as passports to verify the identity. That responsibility is not held by any one authority and is not checked by any agency in the domain key structure.

Another downside is that the use of domain keys has to be a net standard in order to gain any real advantage over spammers. Furthermore spam will not be stopped by this, blacklisting IP-addresses, is almost the same function and has not shown itself to be effective [49].

The up side is that it will probably stop the stupid spammers that cannot bypass a digital identification system, but the smart ones will not be waylaid that easily.

Anecdote

Spam is a social problem therefore it cannot be wholly remedied with a technological solution. Good examples are bank robbers. Despite all the security measures at a bank they still choose to rob them and are sometimes successful.

CHAPTER 5 CONCLUSION

Research into the status quo of the spam world indicates that there are many projects originating from various collaborating fields of expertise, like cryptography, psychology, AI etc.

Spam is not an easy problem to solve; various notable individuals have dedicated a lot of their time trying to come up with a solution. The absence of an absolute solution after years of development indicates that there is more to it than just filtering out messages.

It is wise to assume that what we know theoretically should be tested if possible, to falsify theories and retain facts. One of the problems that arise when dealing with a problem like spam is the social origin. Spam provides a relatively easy income for various people in for example China, Nigeria, etc [52]. It is very popular because it requires a low investment cost with high potential yields and the local government does not do much to actively combat this form of misuse. Leading to the probability that spam will never totally cease, because the success rate, however slim, will provide a temptation for people who's alternatives are not as good [50].

The most widely used and most successful type of solution for stopping spam is located at the recipients' end of mail transmission. The recipient identifies and eliminates spam she receives. This is typically done by applications like Spamsleuth, Spamassassin, CRM114 etc. [51, 34, 3] or a default filter in a mail client.

False positives are the major flaw of spam solutions. To possibly stop spam solutions are needed that can reach the 99.9999% accuracy level when pressed and can reach a 99% level without much trouble. The second condition is to accommodate the normal mail user that cannot invest too much effort and time. Because of the value users attach to their ham and the cost false

positives could bring when the mail is important for business the current solutions do not provide enough "accuracy for time spent value".

With modifications and developments in spam solutions it is theoretically possible to reduce false positives drastically. By improving the accuracy of identification through synergy of interoperating solutions it may become negligible. The use of e-mail mine fielding, inoculation, just in time filtering and possible further enhancements are a sound theoretic basis for this statement. These solutions have been partially tested and seem to work well, giving good hope for the future [30]. Further empirical research is needed to verify theory and to reveal the attainability of these hopes.

Spam has hitherto not been waylaid in any major way and has been growing ever since. Large radical changes in the way spammers operate have therefore not been needed. With new changes at hand in the fight against spam it is plausible that spammers will start to react to this change and improve their spam, creating a more resistant type of spam. The phrase, "necessity is the mother of all inventions" seems appropriate in the future development of spam.

CHAPTER 6 EVALUATION

POSITIVE

The research and analysis of the project gives an objective view of the current spam solutions. When starting the project there were not a lot of projects that had assessed this area. The thesis gives a good high-level view of what is available with distinct categorizations of technology. These enable readers to get a good introduction into the subject and acquire pointers that lead to further sources, for more in-depth information.

The thesis is by large a collection of research performed by notable developers in the anti-spam community. This elevates the reliability of the scientific basis for the thesis. There are no unsubstantiated conclusions or statements; by doing the chance that the thesis contains falsehoods is very slim.

Besides having looked at the theoretical background of every technology, an effort was made to see what is effectively being used in real life situations. According to current scientific beliefs of Epistemology, the empirical side of science cannot be separated from the theoretical. This belief is upheld in the thesis.

Assessments have been made about the viability of multiple solution interactions, to enhance the total spam stopping ability. These were not assessed with a numerical weight scale because of time constraints. But the basis for an empirical examination has been created that can determine in an exact way whether certain collaborative solutions will work.

NEGATIVE

There were a lot of influencing factors from outside the project that interfered and hindered the efficiency of work. This did not decrease the quality of the work that was done, but it did hinder the amount of analysis that could be completed.

The research started out as a development project starting literature survey. Within a few weeks it was apparent that the subject was far to complicate for the allotted time span. To ensure quality of the project, it was changed into a literature study.

REFLECTION

The project needed a literature study to orientate on the current scope of solutions. In hindsight, it might have been even more interesting to examine the mathematical probability of success for a particular solution. By doing this, a foundation is created for the development (or rejection) of one of the proposed techniques.

My preference lies in examining the limits of features variance, when still staying in the humanly understandable range. There are enough probability models that can be used to assess the limits of variance that can be applied features, with in the restrictions of the eye-space rule.

The main thing is to keep the scope very small for a continuation of the research. This is needed because the subject is very complex and easily leads the researcher astray into the many different disciplines that are involved.

SELF-ASSESSMENT:

Quality of research result:	9
Quality of thesis:	9
Complexity of research questions:	10
Relevancy of the project to the subjects of the Master Software Engineering:	8

REFERENCES

- [1] J.Lee, *Internet without Spam, is it possible?*, AP Net Abuse Workshop, http://www.apcauce.org/meetings/030825/proceedings/Jaewoong_Lee.pdf , 2003
- [2] Wikipeda, *Spamming*, <http://en.wikipedia.org/wiki/Spamming> , 2004.
- [3] Bill yerazunis, *crm114 the Controllable Regex Mutilator*, <http://crm114.sourceforge.net/> 2004
- [4] Barry Warsaw, *Anti-Spam Techniques at Python.org* , Spam Conference 2004.
- [6] Open Relay Database, <http://www.ordb.org/>
- [7] CAUCE, *How do you define spam?*, <http://www.cauce.org/about/faq.shtml#how>
- [8] Jon Praed, *"Latest Trends in the Legal Fight Against Spammers"*, spam conference 2004.
- [9] Paul Judge, *The State of the Spam Problem*, Educase review 2003.
- [10] Don M. Blumenthal, Federal Trade Commission Federal Trade Commission Anti-Spam Efforts, February 2004
- [11] Wiki the free Encyclopedia, *SPAM*, <http://en.wikipedia.org/wiki/Spam>
- [12] Fight Spam on the Internet!, *What is spam*, <http://spam.abuse.net/overview/whatisspam.shtml>
- [13] S. Cobb, *The Economics of Spam*, ePrivacy Group, February 2003.
- [14] Brian Sullivan, *The Cost of Spam An ISP Perspective*, October, 2003
- [15] David Harris, *Drowning in sewage*, <http://cauce.org/proceedings.htm> 2003-2004.
- [16] CAUCE, *The Problem*, <http://www.cauce.org/about/problem.shtml>
- [17] ePrivacy Group, *Spam By Numbers*, 2003
- [18] Chris Williams, David Ferris, *The Cost of Spam False Positives*, August 2003, Ferris Analyzer Information Service. Report #385,
- [19] Sullivan T., *The More Things Change: Volatility and Stability in SPAM Features*, MIT SPAM conference 2004
- [20] Spamassassin, *Statistics report for Spamassassin rule set*, <http://spamassassin.apache.org/dist/rules/STATISTICS.txt> 2004
- [21] Gregory L. Wittel, S. Felix Wu, *On Attacking Statistical Spam Filters* July 2004, First Conference on Email and Anti-Spam (CEAS).
- [22] Paul Graham, *Filters vs. Blacklists*, September 2002.
- [23] Sharon Gaudin, *False Positives: Spam's Casualty of War Costing Billion*, Earthweb 2003.
- [24] Philip R. Zimmermann, *Pretty Good Privacy*, <http://www.pgp.com/> 1991.
- [25] Paris Trudeau and Dr. Richard Cullen, Dave Zwieback, *Major Techniques for Classifying Spam*, April 2003
- [26] John Graham-Cumming, *The Spammers' Compendium*, <http://www.jgc.org/tsc/> 2004.

- [27] CAUCE, *Tracking spam*, <http://www.claws-and-paws.com/spam-l/tracking.html>
- [28] Geoff Hulten, Microsoft *Filtering Junk Mail on a Global Scale*, Spam Conference 2004.
- [29] Sharon Gaudin, Suzanne Gaspar, *The Spam police*, Oktober 2001
- [30] William S. Yerazunis, PhD , *The Spam-Filtering Accuracy Plateau at 99.9% Accuracy and How to Get Past It*, 2004.
- [31] Bradley Mitchell, Wireless/Networking, <http://compnetworking.about.com/library/glossary/bldef-hop.htm> 2004.
- [32] Wikipedia The free encyclopedia, *Base64*, <http://en.wikipedia.org/wiki/Base64>
- [33] Miles Libbey, *Learning from 2003: Spamming Trends and Key Insights*, Spam Conference 2004.
- [34] SpamAssassin , <http://spamassassin.apache.org/>
- [35] *Heuristic filters*. http://www.madgoat.com/mx/mx_mgmt_guide.html#heading_8.7
- [36] SPAM archive, <http://www.spamarchive.org/>
- [37] Colin Smith, *Pyzor – How it works*, <http://www.archeus.plus.com/colin/pydoc/overview/> .
- [38] Paul Graham, *A plan for spam*, August 2002.
- [39] Paul Graham, *Better Bayesian filtering*, January 2003
- [40] John Graham-Cumming, “*How to Beat a Bayesian Spam Filter*”, MIT Spam Conference 2004.
- [41] Black spider technologies, *MailControl Spam Technology Overview*, 2004
- [42] R.Gieben, *Chain of trust*, www.miek.nl/publications/thesis/CSI-report.pdf 2003.
- [43] Dspam, *Statistical spam protection*, <http://www.nuclearelephant.com/projects/dspam/> 2003.
- [44] Maureen Francis Mascha, Cathleen L. Miller, *Stop E-mail Snoops*, Journal of accountancy, Technology workshop 2002.
- [45] Eric S. Johansen, Keith Dawson, *Camram*, <http://www.camram.org/camram.pdf> , 2003.
- [46] Adam Back, *Hashcash*, <http://www.hashcash.org/>
- [47] Spamarrest, <http://spamarrest.com/>
- [48] *DomainKeys: Proving and Protecting Email Sender Identity*, Yahoo! Anti spam resource center, <http://antispam.yahoo.com/domainkeys#a9>
- [49] *Yahoo domain keys*, Spam town, <http://spamtown.net/archives/000021.php>
- [50] NetIQ, *Controlling Spam White Paper*, March 2003, http://download.netiq.com/CMS/WHITEPAPER/NetIQ_Controlling%20SPAM%20White%20Paper.pdf
- [51] Spamsleuth, <http://www.bluesquirrel.com/products/SpamSleuth/>
- [52] LI Yuxiao, *Anti-Spam in China*, <http://icauce.org/proceedings.htm> 2004.
- [53] Pyzor, <http://pyzor.sourceforge.net/>
- [54] Susmitha Athota, Vujwala Vangury, Wayne E. Sprague, Alec Yasinsac, *Whitepaper An Overview of Spam Handling Techniques*, <http://taltech.org/SPAM%20Whitepaper.pdf>

[55] Spam in Europe
<http://www.lingsoft.fi/~reriksso/spamcon2001/slide6-0.html> [Era Eriksson]

[57] Email Validation Service, *Ham vs. Spam*,
www.ansci.wsu.edu/help/userhelp/notes/ham_and_spam.pdf

[56] Michael Bevilacqua-Linn, *Machine Learning for Naive Bayesian Spam Filter Tokenization*, 2003,
<http://www.cs.rochester.edu/u/brown/Crypto/studprojs/Spam.pdf>

FIGURE LIST

Figure 1. Spam in Europe 1997-2000	2
Figure 2. General filter error learning curve	11
Figure 3 Pyzor schema [17]	12

APENDIX-A LITERATURE SURVEY FINDINGS

Websites that were found during the survey:

Spammers

1. Astalavista basic security index – <http://www.astalavista.com/>
2. This is a list of the Spam ware vendors and Spam resource suppliers. These are the people who con gullible newcomers into spamming without telling them what the consequences are - http://www.spamsites.org/live_sites.html
3. Known Spammers - <http://sysadmin.info/spamlinks/prospam.htm#spammers>

Anti Spam

4. Microsoft developments against Spam “See Spam run” – <http://research.microsoft.com/displayArticle.aspx?id=411>
5. RSA document “Hot to protect against militant Spammers” - <http://www.rsasecurity.com/rsalabs/staff/bios/mjakobsson/spam/spam.pdf#xml=http://www.rsasecurity.com/programs/texis.exe/webinator/search/xml.txt?query=spam&pr=default&order=r&cq=&id=408b7f542>
6. Everyone.net basic describer of the way Spam is handled by this organisation - <http://www.everyone.net/spam.html>
7. PC world Spam news - <http://www.pcworld.com/resource/spamwatch.asp>
8. UXN Spam combat organisation - <http://combat.uxn.com/>
9. Unicom Software Archive: ungoospam - <http://www.unicom.com/sw/ungoospam/>
10. InternetPrivacy for dummies - <http://www.internetprivacyfordummies.com/modules.php?op=modload&name=Sections&file=index&req=listarticles&secid=3>
11. Anti Spam site, Claws and paws – <http://www.claws-andhttp://www.claws-and-paws.com/spam-l/tracking.html>
12. Anti - Spam abuse site – <http://spam.abuse.net/>
13. Mail Abuse Prevention System LLC and Realtime Blackhole List - <http://www.mail-abuse.org/rbl/>
14. Anti spam info site (make money myths) - <http://www.stopspam.org/>
15. "HOW TO COMPLAIN ABOUT SPAM". An excellent article that talks about various methods to complain about Spam. I highly recommend reading this if you are serious about fighting Spam. - <http://commons.somewhere.com/rre/1997/How.to.Complain.About.Sp.html>
16. Various articles on Spam prevention - <http://internet.designerz.com/internet-abuse-spam-preventing.php>
17. The war on Spam - <http://spam.gunters.org>
18. The Coalition Against Unsolicited E-Mail. CAUCE is a political advocacy group, which is trying to fight Spam on the legal front, as well as keep poorly written bills, which would legitimize Spam from passing. - <http://www.cauce.org>
19. Coalition against unsolicited commercial Email, talking about the definition, what to do when spammed and resources – <http://www.cauce.org/about/faq.shtml>
20. Coalition against unsolicited commercial Email in Europe, talks about specific opportunities to fight Spam in Europe - <http://www.euro.cauce.org/nl/index.html>
21. SPUTUM. The Subgenius Police Usenet Tactical Unit (Mobile) - A spammer's worst nightmare. Their site also has some helpful anti-Spam tutorials and information on persistent Spammers - <http://www.sputum.com/>
22. Expita work in progress trying to provide everything about e-mail as a comprehensive reference tool showing the e-mail user how to make the most of the Internet – <http://www.expita.com>

23. Kill the Spam. Online group fighting Spam - <http://www.studio42.com/kill-the-spam/assistance/>
24. WebSentryTM - http://www.thales-ecurity.com/CaseStudies/Documents/Stampit_Case_Study.pdf
25. A leading figure in the fight against Spam, several papers and archives on Spam - <http://www.paulgraham.com/antispam.html>
26. International Association for Cryptologic Research, a resource for cryptography used against Spam- <http://www.iacr.org/>
27. Spam con foundation - <http://spamcon.org/index.shtml>
28. Spam resource link page - <http://sysadmin.info/spamlinks/spamlinks.htm>

The Spamtools Mailing List

29. A list of responsible Anti-Spam sites which also contains e-mail addresses of the abuse departments for those sites as well as links to their (Anti-Spam) AUPs. Spam boycott a general effort of users to stop Spam. These are users and ISP's etc. They try to stop the use of their resources to spread Spam. IPS also act responsibly towards Spam abuse. - <http://spam.abuse.net/goodsites/>
30. The Net Abuse FAQ with various interesting terminology and history about Spam. - <http://www.faqs.org/faqs/net-abuse-faq/part1/>
31. The E-Mail Abuse FAQ Good definition of what email abuse is. - <http://www.faqs.org/faqs/net-abuse-faq/email-abuse/>

Legal

32. Euro, E-Privacy Directive Proposal COM(2000) 385. - <http://www.euro.cauce.org/en/timeline1.html>

Media

33. Fingerprint, cryptographic identification of mail messages. - http://story.news.yahoo.com/news?tmpl=story&cid=1093&ncid=1093&e=3&u=/pcworld/20040412/tc_pcworld/115640
34. The news section on Yahoo that relates to Spam. Good source for new developments - Yahoo! Full Coverage:Spam Wars - http://headlines.yahoo.com/Full_Coverage/Tech/Spam_Wars/
35. Home of the SPAM News mailing list (daily news and commentary about SPAM) and a collection of anti-SPAM resources. - <http://www.spamnews.com>

Anti-Spam utilities and other technical things.

The first few create fake email fronts for public use, these are filters.

36. AuthentiMail, commercial mail filtrations system - <http://www.despammed.com>
37. A Spam filtering service. Another service which allows you to create "disposable" e-mail addresses. - <http://www.emailias.com/>
38. MailShield is a software plug-in for your existing mail server which can reject Spam, prevent unauthorized mail relaying and halt email bombs. It comes in UNIX and Windows NT flavors. - <http://www.mailshield.com>

These are alternative ways of blocking Spam

39. The Realtime Blackhole List. This is a setup that systems "subscribe" to in order to receive a list of IP addresses and netblocks of sites that either Spam or have open relaying being abused by Spammers which are blocked automatically. - <http://www.mail-abuse.org/rbl>
40. SpamEx. An e-mail service which allows you to create "disposable" e-mail addresses on their system. That way, if you give out such an address to a site that later spams

- you, you can delete that address and stop the flow of Spam from them. - <http://www.spamex.com/>
41. Spam Killer is an e-mail filtering program it also supports methods for writing automatic and manual complaints. - <http://www.spamkiller.com>
 42. Sam Spade. A web based tool which helps novices gather info on spamming domains by performing nslookups, traceroutes, and whois queries on domains and netblocks. - <http://www.samspade.org/>
 43. World Domain Search Registry. A complete list of all top-level domains, with links to the registrar for each domain. - <http://www.uninett.no/navn/domreg.html>
 44. Spam Tracing Resource - <http://sysadmin.info/spamlinks/trace.htm>
 45. SPUTUM tools - <http://www.sputum.com/sputools.html>

English Header Reading Advice

46. The EarthLink Email Protection Agency - <http://www.earthlink.net/epa/spam.html>
47. Spam tracking page - <http://www.rahul.net/falk/index.html>
48. Academic Computing and Communicatios Center - <http://www.uic.edu/depts/accc/newsletter/adn29/headers.html>

Conferences on Spam

49. First Conference on Email and Anti-Spam (CEAS) Mountain View, CA July 30 and 31, 2004 - <http://www.ceas.cc/>
50. "Spam And The Law" Conference Proceedings - <http://www.isipp.org/conference-proceedings.php>
51. Proceedings of the 2003 Spam Conference, papers are available with video streams of the lectures - <http://spamconference.org/proceedings2003.html>
52. Camram is a hybrid antispam system updating the physical world concept of postage to peer-to-peer electronic postage. Camram operates on a peer-to-peer basis using proof of work and digital signature techniques. The camram system intentionally makes payment information visible allowing intermediate machines to filter Spam closer to its ingress. - <http://www.camram.org/>
53. ETC: Email Technology Conference - <http://www.etcevent.com/>

