# Its About You, Me and Every Netizen Because We've Got Spam and Phish!

Shalendra Chhabra
University of California, Riverside
http://www.cs.ucr.edu/~schhabra
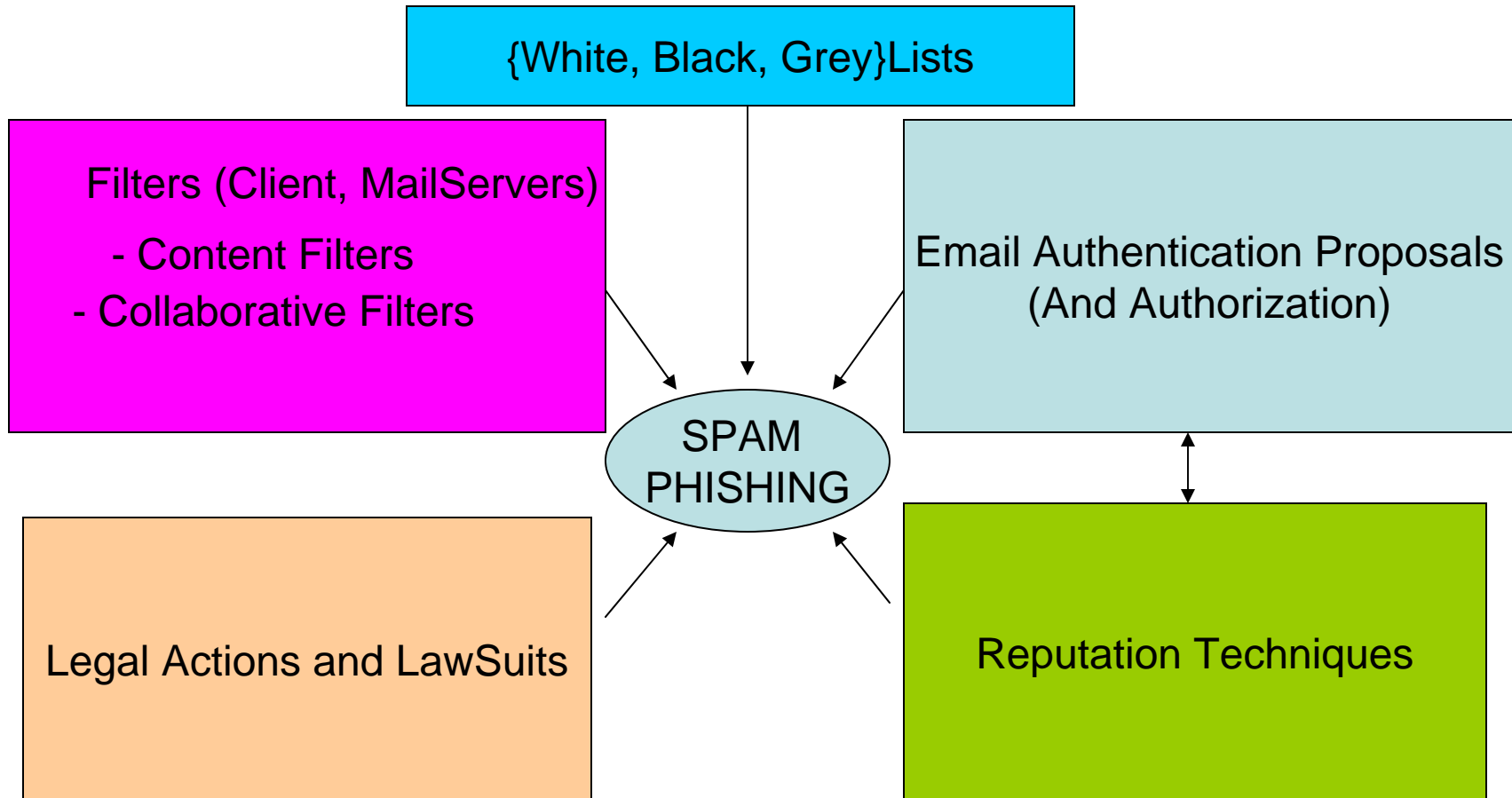http://www.spam-research.com
schhabra@cs.ucr.edu
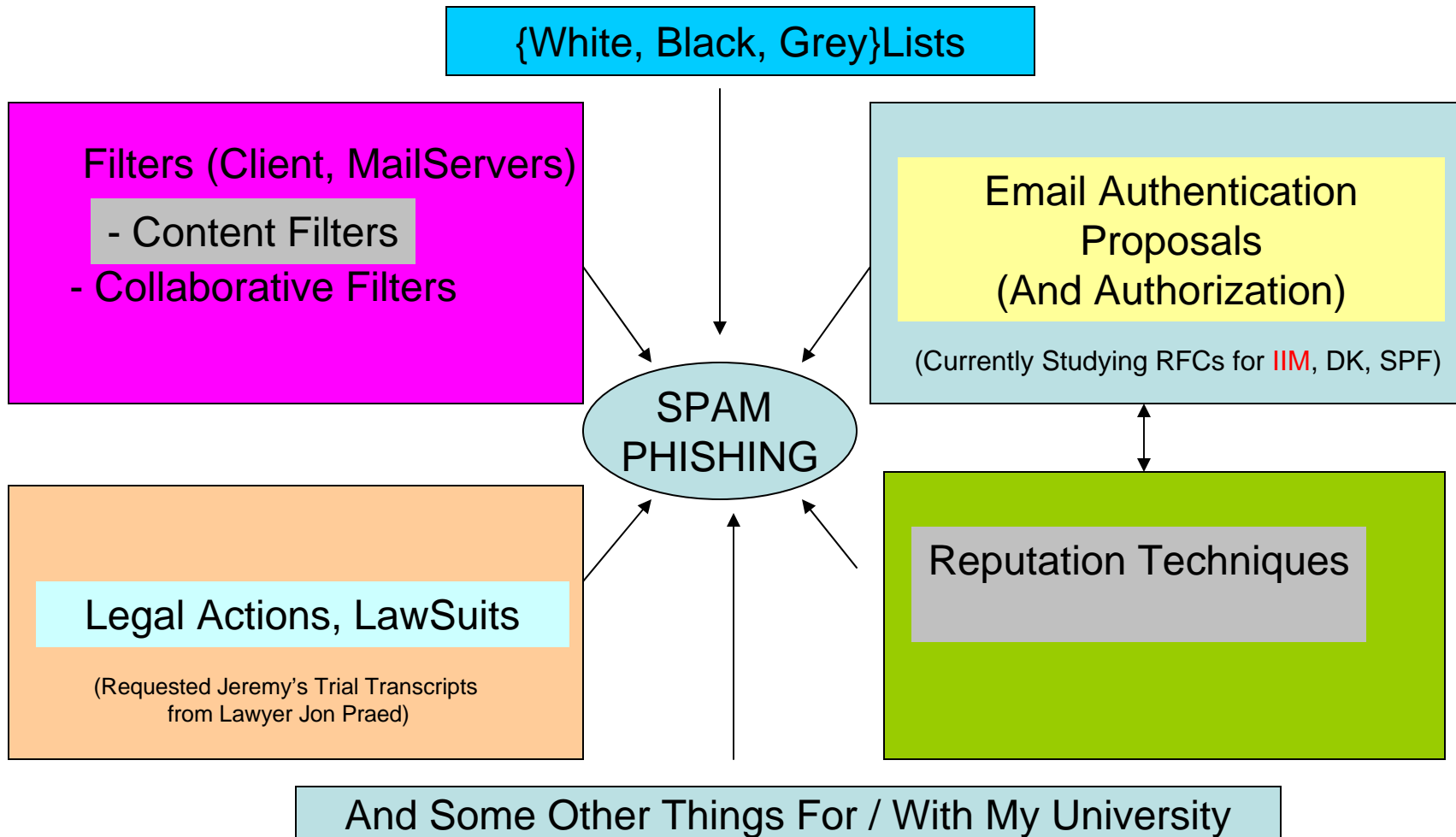
Venue: Cisco Systems

# Its an Honor to Speak Here

- Thanks ☺ to Jim Fenton, Sanjay Pol, Shamim Pirzada and Jennifer Visaya for inviting me

- Regards to Cisco Anti Spam Team Members

- Congratulations to Cisco Systems for acquiring TopSpin

# Tackling Spam and Phishing

{White, Black, Grey}Lists

Filters (Client, MailServers)

- Content Filters
- Collaborative Filters

Email Authentication Proposals
(And Authorization)

SPAM
PHISHING

Legal Actions and LawSuits

Reputation Techniques

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Masters Thesis On Tackling Spam and Phishing

{White, Black, Grey}Lists

Filters (Client, MailServers)
- Content Filters
- Collaborative Filters

Email Authentication
Proposals
(And Authorization)

(Currently Studying RFCs for IIM, DK, SPF)

SPAM
PHISHING

Legal Actions, LawSuits

(Requested Jeremy's Trial Transcripts
from Lawyer Jon Praed)

Reputation Techniques

And Some Other Things For / With My University

04/18/2005

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Motivation and How Did it All Start?

- September 2003 - Once was thinking for a Class Project and got spam, Clicked => Anti Spam
- Heard about MIT Spam Conference, January 2004
- January 2004 - Went up to attend MIT Spam Conference on my own, was a backseat audience
- Spam Conference 2004 - Found some errors in one presentation
- June 2004 - Proposed my Own Model and presented in UK
- 2005 spoke at MIT Spam Conference ☺ on a Unified Model of Spam Filtration

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Bayesian Filters vs Our Model*

- Question: Why not Traditional Pattern Matching Algorithm (KMP) and Suffix Tries ?

- Almost all the filters at MIT Spam Conference Jan 2004, were Naïve Bayesian Filters

- Naïve Bayesian Filters have independence assumption for events for ex:

  *"click here to buy cheap software "* probability of occurrence of *"buy"* is assumed to be independent of probability of occurrence of *"click"* or *"cheap"*

- But probabilities of occurrence of these words together are highly related

- Proposed a Markov Random Field Model where occurrence of one word is dependent on the occurrence of other words in the vicinity, implemented and tested in CRM114

- Accuracy and Performance is higher than Paul Graham's Bayesian Filter Model

*Shalendra Chhabra , William S. Yerazunis, and Christian Siefkes. **"Spam Filtering using a Markov Random Field Model with Variable Weighting Schemas".** In Proceedings of the Fourth IEEE International Conference on Data Mining (ICDM '04), Brighton UK, November 2004.*

Shalendra Chhabra
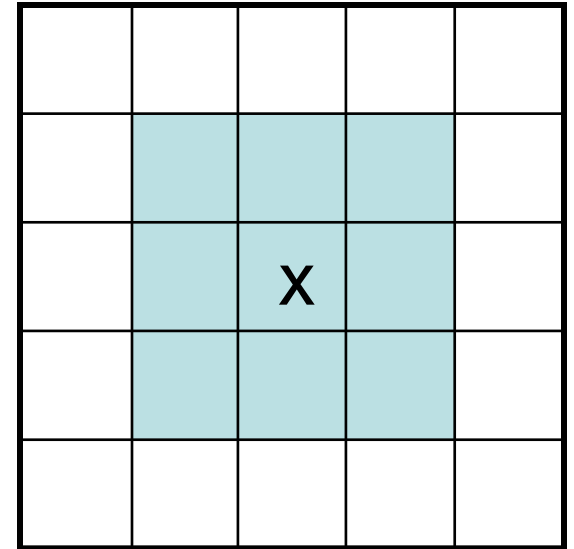(Its About You, Me and Every Netizen -Limited Distribution)

# Borrowed Idea from Computer Vision

- A Site represents a point or region in Euclidean space
- A Label is an event that may happen to a site for ex: In edge detection, the label set is

  L = {edge,non-edge}

- Let $F = \{F_1, F_2, \ldots F_m\}$ be a family of random variables on the discrete set of sites S, in which each random variable $F_i$ takes the value $f_i$ in the discrete label set L

  The family F is called a <u>Random Field</u>

- $P(F = f) = P(F_1 = f_1, F_2 = f_2, F_3 = f_3 \ldots, F_m = f_m)$ denotes a joint event
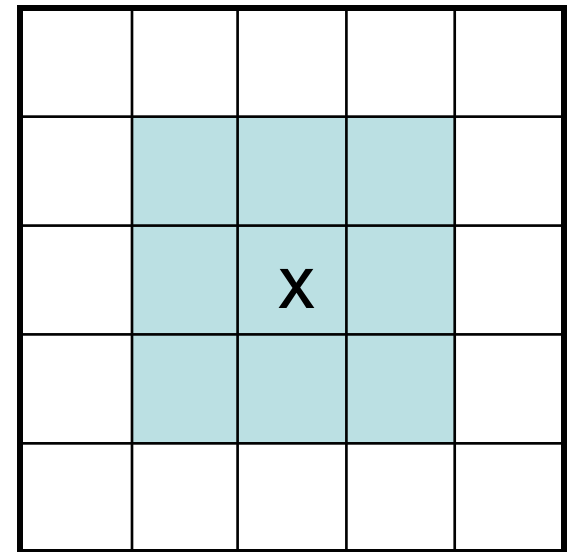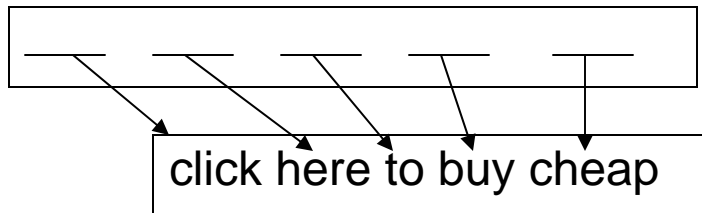
# Neighborhood System

- The Sites in S are related to one another via a Neighborhood System. A Neighborhood System for a site X denotes the set of sites surrounding X

-  Any F is said to be a MRF on S with respect to a neighborhood N iff:

*1. $P(f) > 0$ ;  (positivity)*
*2. $P(f_i|f_{S-\{i\}}) = P(f_i|f_{N_i})$ (Markovianity)*

X

# Analogy with Spam Text

A Site in the context of spam classification refers to
*relative position* of word in a sequence
And a Label maps to *word values*

click here to buy cheap

X

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Assigning Weights to These Features

- Sequence ABC has 8 subsequences including empty sequence and itself:

  {A, B, C, A_C, BC, AB, ABC, 0}.

- Idea: Weight of Feature with n terms in the sequence should be greater than combined weight of all Features of length less than n:

$$W(n) > \sum_{k=1}^{n-1} \left( \binom{n}{k} \times W(k) \right)$$

# Weighting Schemes

### Minimum Weighting Schemes

### Exponential Scheme

$$W(n) = \sum_{k=1}^{n-1} \left( \binom{n}{k} \times W(k) \right) + 1.$$

$$base^{n-1} > \sum_{k=1}^{n-1} \left( \binom{n}{k} \times base^{k-1} \right)$$

| n | MWS | ES |
|---|---|---|
| 1 | 1 | 1 |
| 2 | 1, 3 | 1, 3 |
| 3 | 1, 3, 13 | 1, 5, 25 |
| 4 | 1, 3, 13, 75 | 1, 6, 36, 216 |
| 5 | 1, 3, 13, 75, 541 | 1, 7, 49, 343, 2401 |
| 6 | 1, 3, 13, 75, 541, 4683 | 1, 8, 64, 512, 4096, 32768 |

**Table 1. Minimum & Exponential Weightings**

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Example Subphrases and Models Tested

| n | MWS | ES |
|---|---|---|
| 1 | 1 | 1 |
| 2 | 1, 3 | 1, 3 |
| 3 | 1, 3, 13 | 1, 5, 25 |
| 4 | 1, 3, 13, 75 | 1, 6, 36, 216 |
| 5 | 1, 3, 13, 75, 541 | 1, 7, 49, 343, 2401 |
| 6 | 1, 3, 13, 75, 541, 4683 | 1, 8, 64, 512, 4096, 32768 |

**Table 1. Minimum & Exponential Weightings**

| Text | SBPH | ESM | MWS | ES |
|---|---|---|---|---|
| Do | 1 | 1 | 1 | 1 |
| Do you | 1 | 4 | 3 | 8 |
| Do $<skip>$feel | 1 | 4 | 3 | 8 |
| Do you feel | 1 | 16 | 13 | 64 |
| Do $<skip><skip>$lucky? | 1 | 4 | 3 | 8 |
| Do you $<skip>$lucky? | 1 | 16 | 13 | 64 |
| Do $<skip>$feel lucky? | 1 | 16 | 13 | 64 |
| Do you feel lucky? | 1 | 64 | 75 | 512 |

SBPH:        1,1,1,1,1

ESM ($2^{2(n-1)}$): 1,4,16,64

# MRF Model for Spam

- All incoming email is broken in features
- *A random class function C is defined C:Omega -> {spam,nonspam}*

- $P(spam|F_i) = P(F_i|spam)P(spam)$
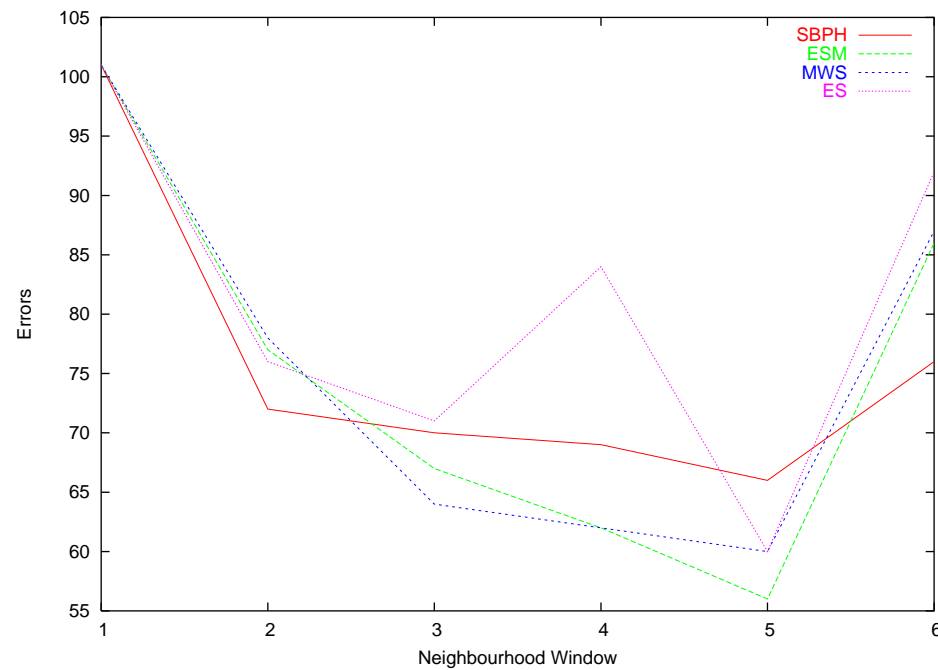
$$-------------------------------$$

$(P(F_i |spam)P(spam)+P( F_i|ham)P(ham) )$

- *Local Formula for $P(F_i|spam)$ \**

- *The output $P(spam|F_i)$ becomes $P(spam)$ for the feature $F_{i+1}$*

  <span style="color:red">If $P(spam|F_n)$ is higher than $P(ham|F_n)$ , email is tagged as "spam"</span>

# Results with MRF Model for Spam Filtering

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

## Winnow Algorithm and Orthogonal Sparse Bigrams**

- Winnow is a statistical but non probabilistic algorithm i.e. it computes score and not probability
- It keeps n dimensional weight vector for each class c, i.e. $w^c=(w^c_{1,} w^c_2, \ldots w^c_m)$, where $w^c_i$ is the weight of the $i^{th}$ feature for class c
- The algorithm returns 1 for a class iff the summed weights for all active features surpass a predefined threshold

*** Christian Siefkes, Fidelis Assis, Shalendra Chhabra and William S. Yerazunis.* **Combining Winnow and Orthogonal Sparse Bigrams for Incremental Spam Filtering.** *Lecture Notes in Computer Science. Springer, 2004*, Springer Verlag

# Expressivity of Features

**Table 2.** Features Generated by SBPH and OSB

| Number | SBPH | | | | | OSB | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 (1) | | | | | today? | | | | | |
| 3 (11) | | | | lucky | today? | | | | lucky | today? |
| 5 (101) | | | feel | *<skip>* | today? | | | feel | *<skip>* | today? |
| 7 (111) | | | feel | lucky | today? | | | | | |
| 9 (1001) | | you | *<skip>* | *<skip>* | today? | | you | *<skip>* | *<skip>* | today? |
| 11 (1011) | | you | *<skip>* | lucky | today? | | | | | |
| 13 (1101) | | you | feel | *<skip>* | today? | | | | | |
| 15 (1111) | | you | feel | lucky | today? | | | | | |
| 17 (10001) | Do | *<skip>* | *<skip>* | *<skip>* | today? | Do | *<skip>* | *<skip>* | *<skip>* | today? |
| 19 (10011) | Do | *<skip>* | *<skip>* | lucky | today? | | | | | |
| 21 (10101) | Do | *<skip>* | feel | *<skip>* | today? | | | | | |
| 23 (10111) | Do | *<skip>* | feel | lucky | today? | | | | | |
| 25 (11001) | Do | you | *<skip>* | *<skip>* | today? | | | | | |
| 27 (11011) | Do | you | *<skip>* | lucky | today? | | | | | |
| 29 (11101) | Do | you | feel | *<skip>* | today? | | | | | |
| 31 (11111) | Do | you | feel | lucky | today? | | | | | |

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Comparison of Winnow, Naïve Bayes and CRM114 MRF Model

| Store Size | Naive Bayes<br>All | CRM114<br>$1048577\ (2^{20}+1)$ | CRM114<br>All | Winnow+OSB<br>All |
|---|---|---|---|---|
| Last 500 | 1.84% (9.2) | 1.12% (5.6) | 1.16% (5.8) | **0.46% (2.3)** |
| All | 3.44% (142.8) | 2.71% (112.5) | 2.73% (113.2) | **1.30% (53.9)** |

Note that Error Rate is Halved and Computational Overhead is also reduced
(retaining the expressivity)

# A Unified Model of Spam Filtration
# MIT Spam Conference, 2005

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Pre Processing: Arbitrary Text to Text Transformation

➢ Character Set Folding / Case Folding

➢ Stopword Removal

➢ MIME Normalization / Base64 Decoding

➢ HTML Decommenting

  Hypertextus Interruptus

➢ Heuristic Tagging

  "FORGED_OUTLOOK_TAGS"

➢ Identifying Lookalike Transformations

  '@' instead of 'a', $ instead of S

  Ex: V1agra

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# A Unified Model of Spam Filtration
# MIT Spam Conference, 2005

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Tackling Spam and Phishing

{White, Black, Grey}Lists

Filters (Client, MailServers)

- Content Filters
- Collaborative Filters

**Email Authentication Proposals
(And Authorization)**

SPAM
PHISHING

Legal Actions and LawSuits

Reputation Techniques

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Authentication and Authorization

- Authentication is the process of checking or verifying an entity using some form of integrity information such as an authorization policy.

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Cisco's IIM

**Sending Domain checks if the Source is allowed to send Mail using its Domain**

Analysis of Reputation Attacks (Adapt *IDStealth, Shilling, PseudoSpoofing* and Check )

Typical Identified Internet Mail Message Flow

② Receiver's MTA verifies message

① Sender's MTA signs message

**Sending Domain**

③ Key Registration Server (or DNS) returns result

④ Optional: Consult 3rd party reputation service

**Receiving Domain**

HTTP Query: PlainText
HTTPS – SSL, TLS?
Response Format with values not mentioned in RFC
(Locally Sensitive Hash) ex: Nilsimsa Hash?

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# State with Email Authentication Systems *
## (John Graham Cumming)

incoming email

Has sender
authentication
record? — no → N/A

yes

Sender
authentication
record
matches? — no → negative

yes

positive

**Forged Message or False Negative**

Use Bayesian Filter to Train (State, Output) ☺

**Only sure when its positive: like whitelists**

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# With Email Authentication Systems What's Going to Happen Next?

- Spammers are adept at deploying sender authentication technologies for domains they are not forging

- Timeliness /reputation of domain in place

- Spammers will send from non-forged addresses (Blacklisting is the solution)

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

**Initiator p**  ...  **UltraPeers U, servents s**

Query(*search_string,min_speed*)

$p \longrightarrow *$

QueryHit(*num_hits,IP,port,speed,Result,servent_id$_i$*)

$s_i \longrightarrow p, \ (\forall s_i \in O)$

Select top list $T$ of offerers
Generate a pair ($PK_{poll}$,$SK_{poll}$)

$p \longrightarrow *$

PollRequest(*T*,$PK_{poll}$)

CumulativePollReply({$(PK_i, IP, port, votes, PK_{poll}, servent\_id_i, sgn)$}$_{PK_{poll},...}$)

$v_i \longrightarrow p, \ (\forall v_i \in U)$

(a)

Select a random set $V'$ from the elected voters/clusters
Remove suspicious voters from set $V$

$p \xrightarrow{D} v_j, \ (\forall v_j \in V')$

TrueVote(*Votes$_j$*)

TrueVoteReply(*response*)

$v_j \xrightarrow{D} p, \ (\forall v_j \in V')$

If *response* is negative, discard *Votes$_j$*

Based on valid votes select servent $s$ from which download files

(b)

$p \longrightarrow v_j, \ (\forall v_j \in V')$

PushVote(*servent_id, IP, port,*{TrueVote}$_{PK_{v_j}}$)

$v_j \xrightarrow{D} p, \ (\forall v_j \in V')$

TrueVoteReply(*response*)

(c)

$p \longrightarrow v_j, \ (\forall v_j \in V')$

PushVote(*servent_id, (IP, port)$_r$,*TrueVote)

repeater r

TrueVoteReply(($IP, Port$)$_{PK_{v_j}}$,*response*)

$v_j \xrightarrow{D} r, \ (\forall v_j \in V')$

connection

TrueVoteReply(($IP, Port$)$_{PK_{v_j}}$, *response*)

(d)

**Initiator p**  ...  **Servent s**

Generate a random string $r$

$p \xrightarrow{D} s$

challenge(*r*)

response([$r$]$_{SK_s}$,$PK_s$)

$s \xrightarrow{D} p$

If h($PK_s$)=*servent_id$_s$* $\wedge$ {[$r$]$_{SK_s}$}$_{PK_s}$ = $r$: download
Update experience_repository

(e)

**Figure 1.** SupRep protocol: query and poll (a), vote verification (b)-(d), and resource download (e)

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Check Possibility of These Attacks when using Third Party Reputation Services with Email Authentication Systems

- *PseduoSpoofing*: Forging great number of votes from a single node, giving them different IP addresses, and multiple IDs (TrueVoteConnection detects this)
- *Shilling*: Clique / Control over many servents affecting reputation (Scalability in Gnutella and repeaters for servents behind firewalls takes care of this)
- *ID Stealth*: Malicious Servent replies with QueryReplie's as if generated from genuine servents (Challenge Response detects this)

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Lessons from the Past

- Always think about the possibility of DNS Poisoning in Caches    (Refer **Using** the **Domain Name System** for **System** Break-ins - Bellovin)

- IP Spoofing Attacks

- DoS Attacks on Blacklists

- Some other Ideas ex: LOC record in DNS (Zombie Zones)

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Other stuff I am doing

- Conducting a survey at UCR ( population > 10000 ) – This will give us an idea how students and professors react to spam (will publish in *Nature*)

- Implementing Spam Filters at UCR MailServers in cooperation with the author of these filters and write effective guidelines for system administrators

- antispam.ucr.edu , antispam.cs.ucr.edu

- Yahoo Mail SpamGuard SplitFit   Yahoo!
  (with Miles Libbey)

- A comment on Microsoft's Article on Slashdot
  (On Nilsimsha Hash and "Cmabirgde Uinersvtiy Sapm". , It was on Slashdot )

# On Slashdot

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Spam-Research

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Finishing My Thesis

- Want to make my thesis a very important resource for Anti Spam Industry
- And Miles to go before I sleep….

  In order to contribute have to learn a lot with disciplined and ambitious instincts

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Seek Your Blessings, Guidance, Comments and Criticism for becoming an Anti Spam Leader within next 5 years

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)

# Spam Free World?



Thank You!

Shalendra Chhabra
(Its About You, Me and Every Netizen -Limited Distribution)