



Detecting Sensitive Data Exfiltration by an Insider Attack

Dipak Ghosal
University of California, Davis

Collaborators

- Tracy Liu (PhD Student, UCDavis)
- Rennie Archibald (PhD Student, UCDavis)
- Matt Masuda (Undergraduate Student, UC Davis)
- Cherita Corbett (Sandia National Labs – Livermore)
- Ken Chiang (Sandia National Labs – Livermore)
- Raj Savoor (AT&T Labs)
- Zhi Li (AT&T Labs)
- Sam Ou (ex AT&T Labs)

6/17/08

NSF I/UCRC



Outline

- Application Identification
- Content Signature Generation and Detection
- Detecting Covert Communication
- Research Directions

6/17/08

NSF I/UCRC



Insider Attack and Insider Threat

- Insider attack
 - *"The potential damage to the interests of an organization by a person who is regarded, falsely, as loyally working for or on behalf of the organization, or who inadvertently commits security breaches."*
- An insider attack can occur through
 - Inadvertent security breach by an authorized user
 - A planned security breach by an authorized user
 - A compromised system by an outsider

6/17/08

NSF I/UCRC



Signals

- **Inter-arrival time:** derived from the sequence of timestamps noted by the sniffer for packets inbound to the host
- **Inter-departure time:** derived from the sequence of timestamps noted by the sniffer for packets outbound from the host
- **Incoming packet size:** vector of packet sizes for HTTP packets inbound to the host
- **Outgoing packet size:** vector of packet sizes for packets outbound from the host
- **Outgoing Discrete Time Total Bytes:** vector of outgoing bytes of data aggregated over discrete and fixed time bins

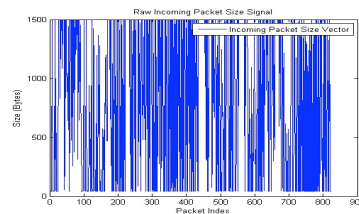
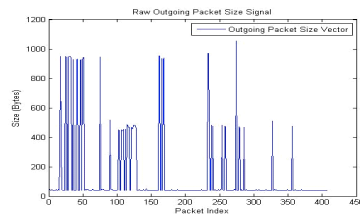


6/17/08

NSF I/UCRC

Signals – Examples

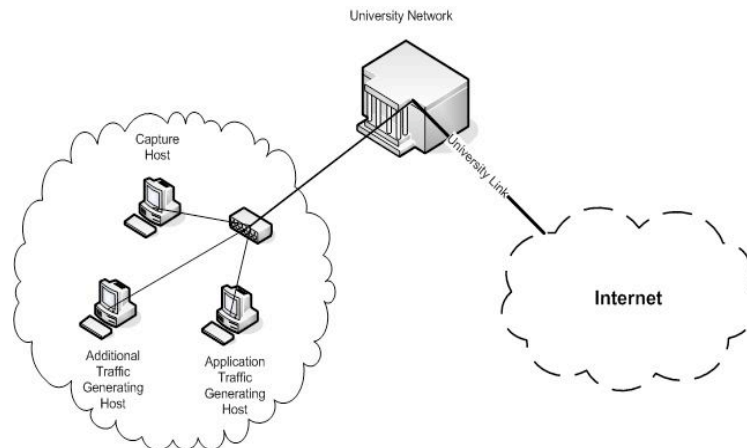
- Outgoing packet size vs. incoming packet size



6/17/08

NSF I/UCRC

Experimental Setup

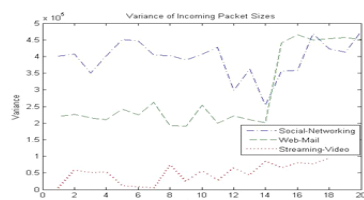


6/17/08

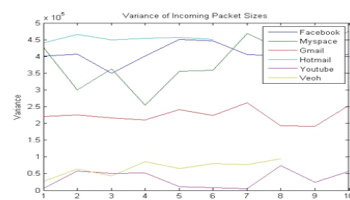
NSF I/UCRC

9

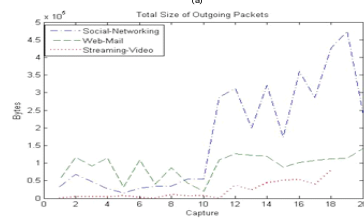
Temporal Statistics



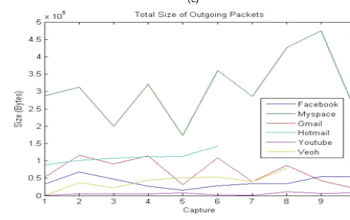
(a)



(e)



(b)



(d)

6/17/08

NSF I/UCRC



Temporal Characteristics

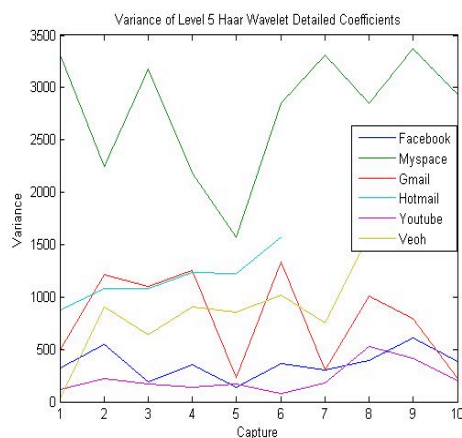
	Interdeparture		Outgoing Packet Size				Incoming Packet Size			
	Q3		Mean		Var		Var		Total Bytes	
	Mean ($\times 10^{-2}$)	Var ($\times 10^{-5}$)	Mean ($\times 10^2$)	Var ($\times 10^3$)	Mean ($\times 10^4$)	Var ($\times 10^9$)	Mean ($\times 10^4$)	Var ($\times 10^7$)	Mean ($\times 10^5$)	Var ($\times 10^{11}$)
Facebook	6.04	38.1	1.15	1.87	5.59	1.73	39.9	119	9.55	2.33
Myspace	1.94	3.1	3.27	8.37	28.6	4.89	38.4	509	31.1	29.9
Gmail	21.54	650	2.71	1.41	18.8	16.7	21.9	46.1	8.73	2.1
Hotmail	6.2	17.2	2.18	0.598	18	4.29	45.3	7.14	7.19	0.432
Youtube	11.26	250	0.82	1.35	4.53	14.9	3.47	69.4	121	133
Veoh	5.43	54.5	0.96	0.75	4	4.17	6.71	49.5	169	118

6/17/08

NSF I/UCRC

11

Wavelet Analysis



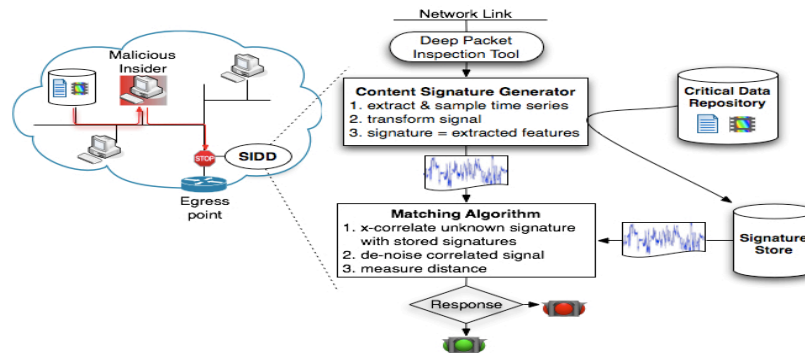
- Use Haar wavelet
- Feature used for comparison
 - Variance of the Level-5 detailed coefficients

6/17/08

NSF I/UCRC

12

Content Identification: Motivation



Can we detect illegal dissemination of protected digital (media) assets?

6/17/08

NSF I/UCRC

13

Content Signature

■ Content-based Signature

- “The media itself is a watermark”
- Unique and robust
 - Different content should have distinct signatures
 - The signatures are tolerant to various forms of noise and distortions
- Requirements vary with applications
 - From video search to detecting video copying

6/17/08

NSF I/UCRC



Content Signature Generation

■ Basic idea

- Extract a time series (or signal) of the content and analyze the signal to generate the signatures
- Capture the temporal correlation in the signature
- Treating the content signatures as time series
 - Use signal processing techniques and tools to analyze
 - Wavelet transform
 - Any portion of the content can be used for detection
 - Computation cost saving

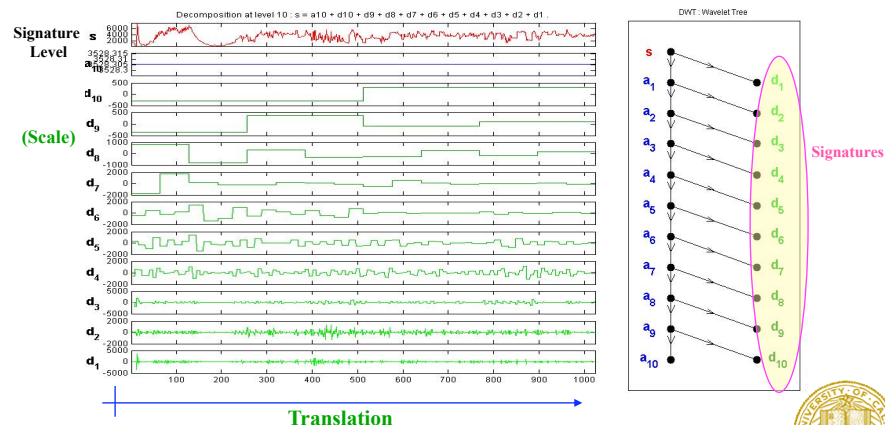


6/17/08

NSF I/UCRC

Content Signature Generation – Example

■ The Detailed Coefficients of the Star Wars Movie



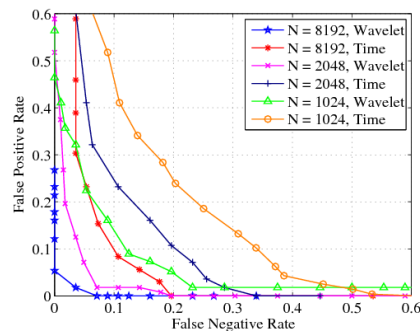
6/17/08

NSF I/UCRC

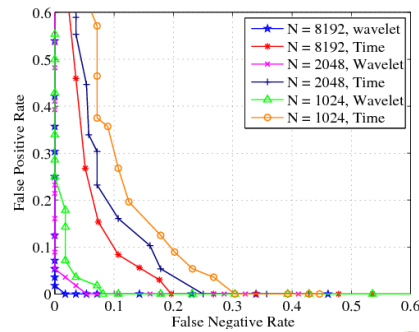


Preliminary Analysis

ROC curve in rate adaption case 1



ROC curve in rate adaption case 2



6/17/08

NSF I/UCRC



Detecting Covert Communication

- Exfiltration of sensitive information may be carried out using covert communication
 - Hiding content/communication in an innocuous carrier using a steganography tool
- Challenges
 - The content may be encrypted
 - Different types of carriers

6/17/08

NSF I/UCRC



Audio Steganalysis

- The analysis and classification method of determining if an audio bears hidden information
- Easy to establish
 - Voice over Internet Protocol (VoIP) and other Peer-to-Peer (P2P) audio service
- High hidden capacity
 - Inherent redundancy in the audio signal
 - Its transient and unpredictable characteristics
- Human ear is insensitive to small distortions

6/17/08

NSF I/UCRC



Main Points

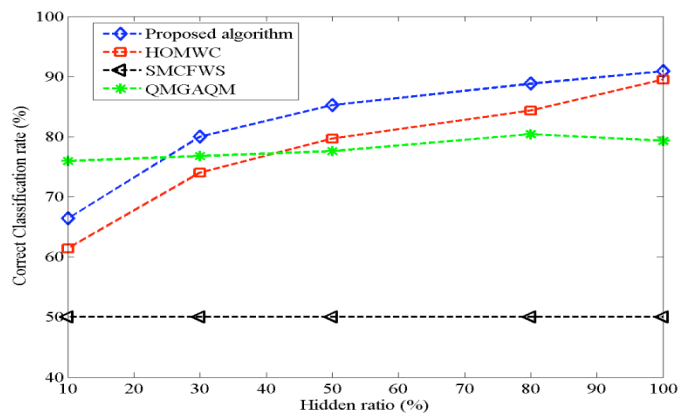
- A new approach to detect hidden content in audio files
- Uses Hausdorff distance and feature vectors based on higher-order statistics
- Good detection rate even with low hidden ratio

6/17/08

NSF I/UCRC



Comparative Analysis



6/17/08

NSF I/UCRC



Research Directions

- Improving the techniques
 - Wavelet analysis allows time frequency localization
 - Where approximately time certain frequencies occur
 - Is it useful in disambiguating applications?
 - Co-integration can extract similarities in signals that may be uncorrelated
 - Can this be used to detect content that is encrypted and/or modified to evade detection?
- Developing prototypes
 - A VoIP steganalysis tool
 - A classifier for network level application identification

6/17/08

NSF I/UCRC

