# CE 817 - Advanced Network Security
# Network Forensics

Lecture 22

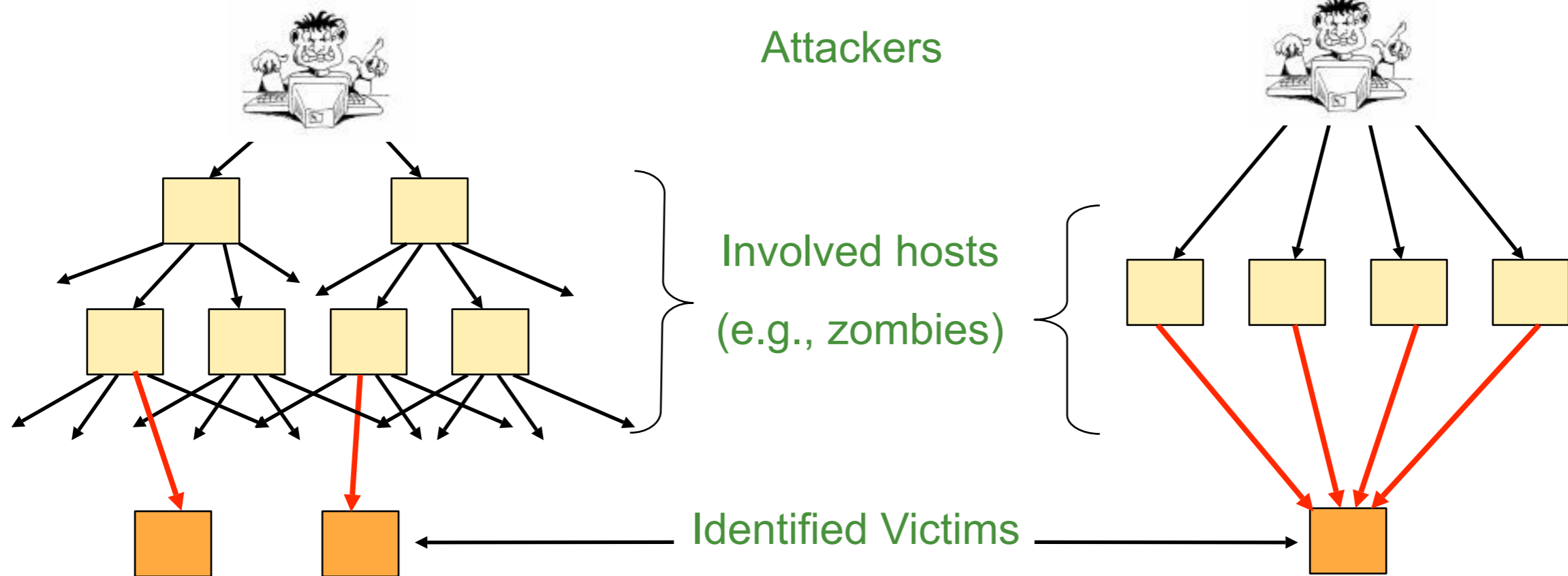Mehdi Kharrazi
Department of Computer Engineering
Sharif University of Technology

# The Structure of Attacks



Attackers

Involved hosts

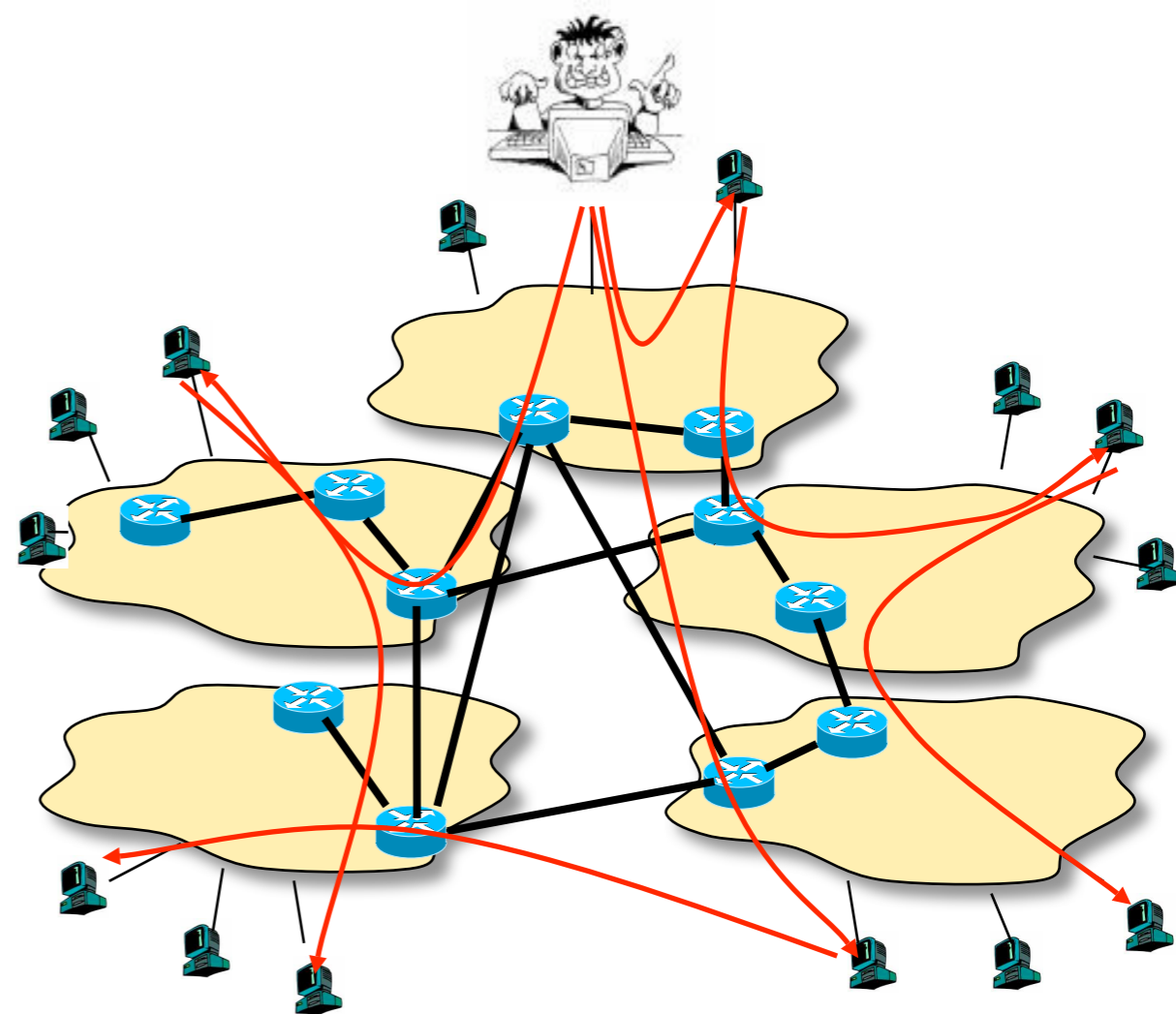(e.g., zombies)

Identified Victims

**Worm Infection**

**Distributed DoS**

- Modern attacks are multi-level
  - Large scale: difficult to defend
  - Hidden trail: difficult to identify initial launch point

# Existing Approaches

- Firewall

- Intrusion detection

  - Identify compromised hosts
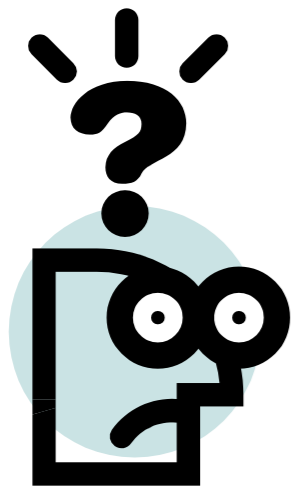
- IP traceback

  - Trace to the source of a packet

# Question:

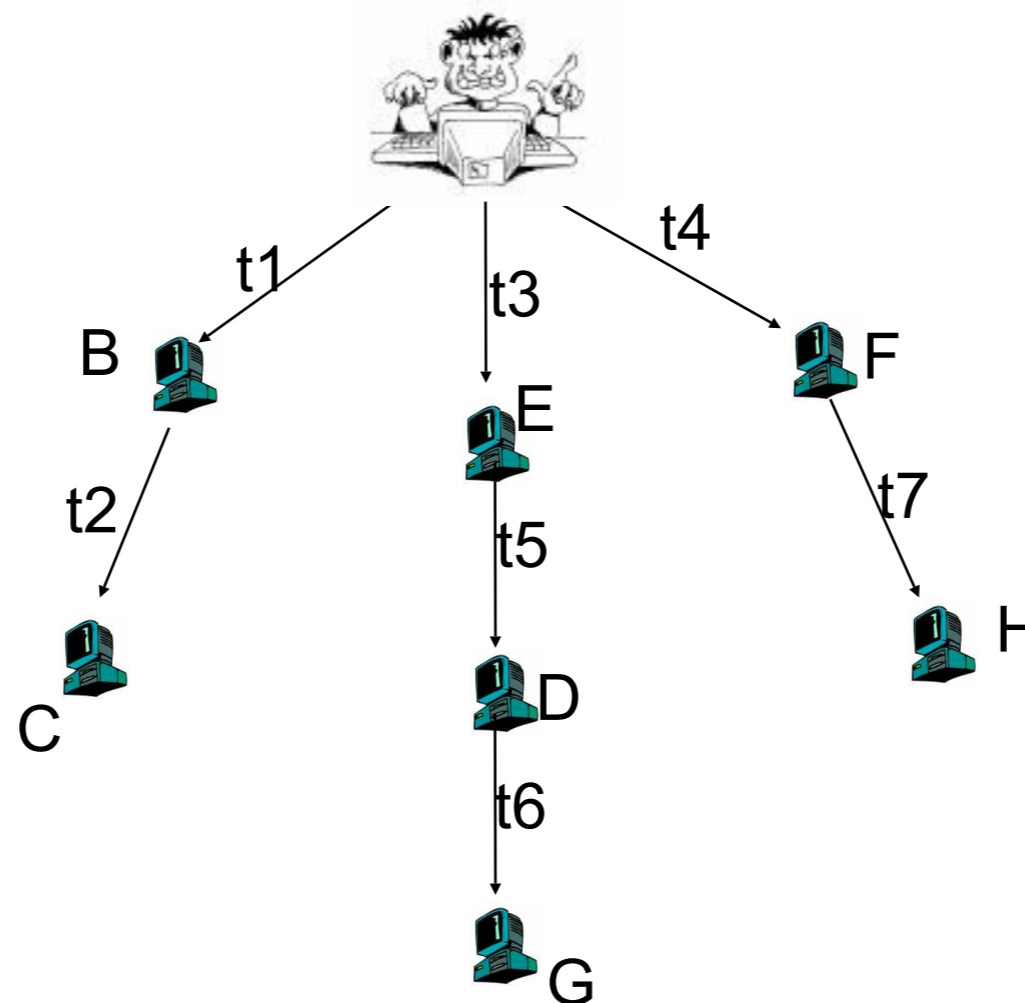- What fundamental capabilities can be added to the network to achieve better security?

# Forensics

- A fundamental capability --- network auditing and forensic analysis
  - Keep communication records
  - Permit post-mortem analysis of patterns across network and time
- Scope: Internet and intranet
  - Correct weak points in a network perimeter
  - Deter future similar attacks

# Two Applications



- Attack Reconstruction: infer which communication carry the attack forward
- Attacker Identification: pinpoint the attack source(s)

# Key Question 1: Feasibility of Network Auditing

- How much storage is needed for auditing?
  - Directional network flows
  - <source, destination, start-time, end-time>
  - A large ISP with O(100) POPs – 450 GB/hour with compression
  - Intranet requirement will be smaller

# Key Question 2: Feasibility of Network Forensic

- Given complete information, can we find needles in the haystack?
  - Host contact graph is large and noisy
  - Algorithms to identify global correlations

# Payload Attribution via Hierarchical Bloom Filters,

Kulesh Shanmugasundaram, Hervé Brönnimann, and Nasir Memon. ACM Computer Communications and Security (CCS'04), Washington, DC, 2004.

# Payload Attribution

- The problem:

    - Identify the sources and/or the destinations of a bit-string in a network

    - We may only have an arbitrary portion of payload

# Challenges Facing Network Forensics

- Lack of Infrastructure:

  - For data collection, archival, and dissemination

- Volume of Data:

  - Prolonged storage, processing, and sharing of raw data infeasible

  - Even a network of ~3000 hosts have a ~1TB/day requirement!

- Process is Manual:

  - Spans multiple administrative domains

  - Response times very long (digital evidence disappears)

- Unreliable Logging Mechanisms:

  - Host logs are usually compromised

  - Growing support for mobility makes it difficult to maintain prudent logging policies on hosts
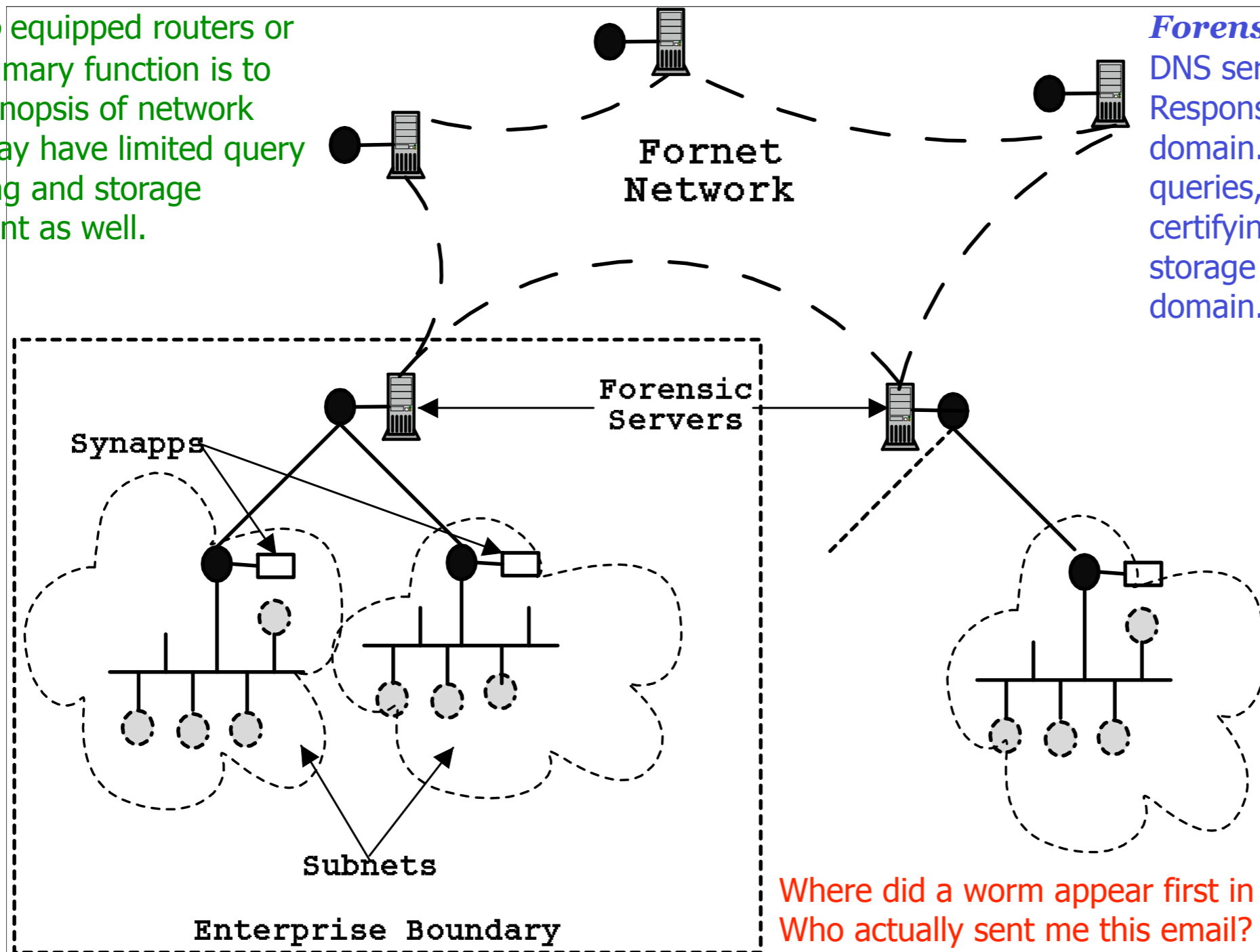
# Our Solution…

- Securely collect, store, disseminate, and process synopsis of network traffic.

- In other words a device analogous to a surveillance camera for the network. Even better – a co-operating network of such surveillance cameras.

- Goal of Project ForNet: development of tools, techniques, and infrastructure to aid rapid investigation and identification of cyber crimes

# ForNet Vision - A unified, end-to-end approach for network monitoring

**SynApp** equipped routers or hosts. Primary function is to create synopsis of network traffic. May have limited query processing and storage component as well.

Fornet Network

**Forensic Server:** Like the DNS servers' of today. Responsible for a network domain. Responsible for routing queries, query processing, certifying results, and long-term storage synopsis for the domain.

Synapps

Forensic Servers

Subnets

Enterprise Boundary

Where did a worm appear first in the .edu domain?
Who actually sent me this email?
When and how was this file transferred to this host?
Was this host compromised in the past? When and how?

# What is a Synopsis?

Data structures and algorithms for representing a set of elements succinctly with predefined loss in information and has the ability to recall the original set of elements with a preset accuracy.

# Properties of a Good Synopsis

- Contains enough data to answer certain classes of queries

  - Who sent payload "xyz"?

  - What did host bug.poly.edu send?

- Contains enough data to quantify confidence of its answers

  - I'm 99.37% sure bug.poly.edu sent "XYZ"

- Have small memory footprint and easy to update

  - Need 20GB/day to keep 1TB/day of raw network data

  - Need to compute a number of hashes per packet

- Resource requirements are tunable

  - Can only afford 3GB/day, adjust the accuracy to accommodate this.

# Advantages of Using Synopses
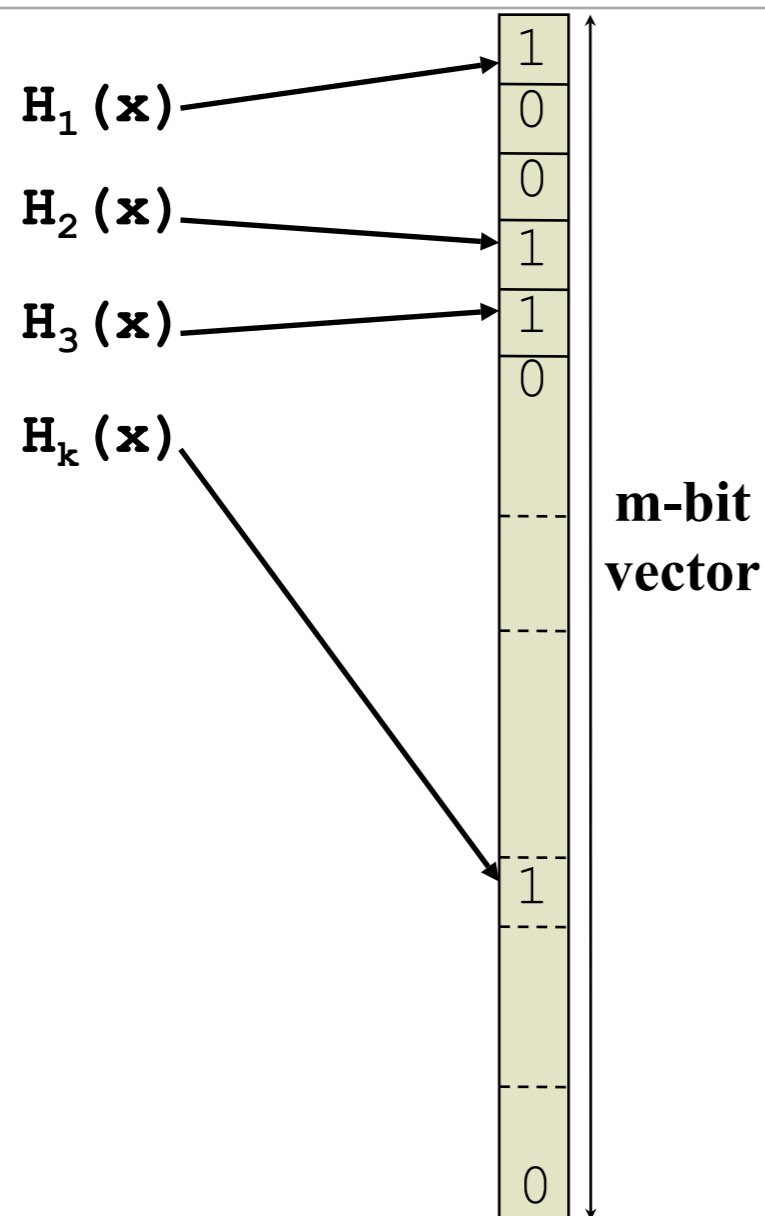
> ## Can retain potential evidence for months!

- Succinct representation of raw data makes it possible to transfer network data to disks

- Sharing/transferring raw data over network is impossible but synopsis can be moved to remote sites

- Query processing would be expensive with raw data

  - What's the frequency of traffic to port 80 in the past week? (raw data vs. a histogram)

- Easily adaptable to various resource requirements

  - Can adopt the size, processing requirements of a Bloom Filter based on various hardware resources and network load

# Bloom Filters

$H_1(\mathbf{x})$

$H_2(\mathbf{x})$

$H_3(\mathbf{x})$

$H_k(\mathbf{x})$

| |
|---|
| 1 |
| 0 |
| 0 |
| 1 |
| 1 |
| 0 |

**m-bit vector**

| |
|---|
| 1 |

| |
|---|
| 0 |

$$FP = \left( 1 - (1 - 1/m)^{kn} \right)^k$$

- Bloom Filter:
  - Randomized data structure for representing a set in order to support membership queries.
  - Insert(x):
    - Flip bits H1(x)… Hk(x) to '1'
  - IsMember(y):
    - If H1(Y) … Hk(Y) all '1' "yes" otherwise "no"
- Can tradeoff memory (m), compute power (k), and accuracy (FP)
  - m – length of bit vector (range of H(.))
  - k – number of hashes per element
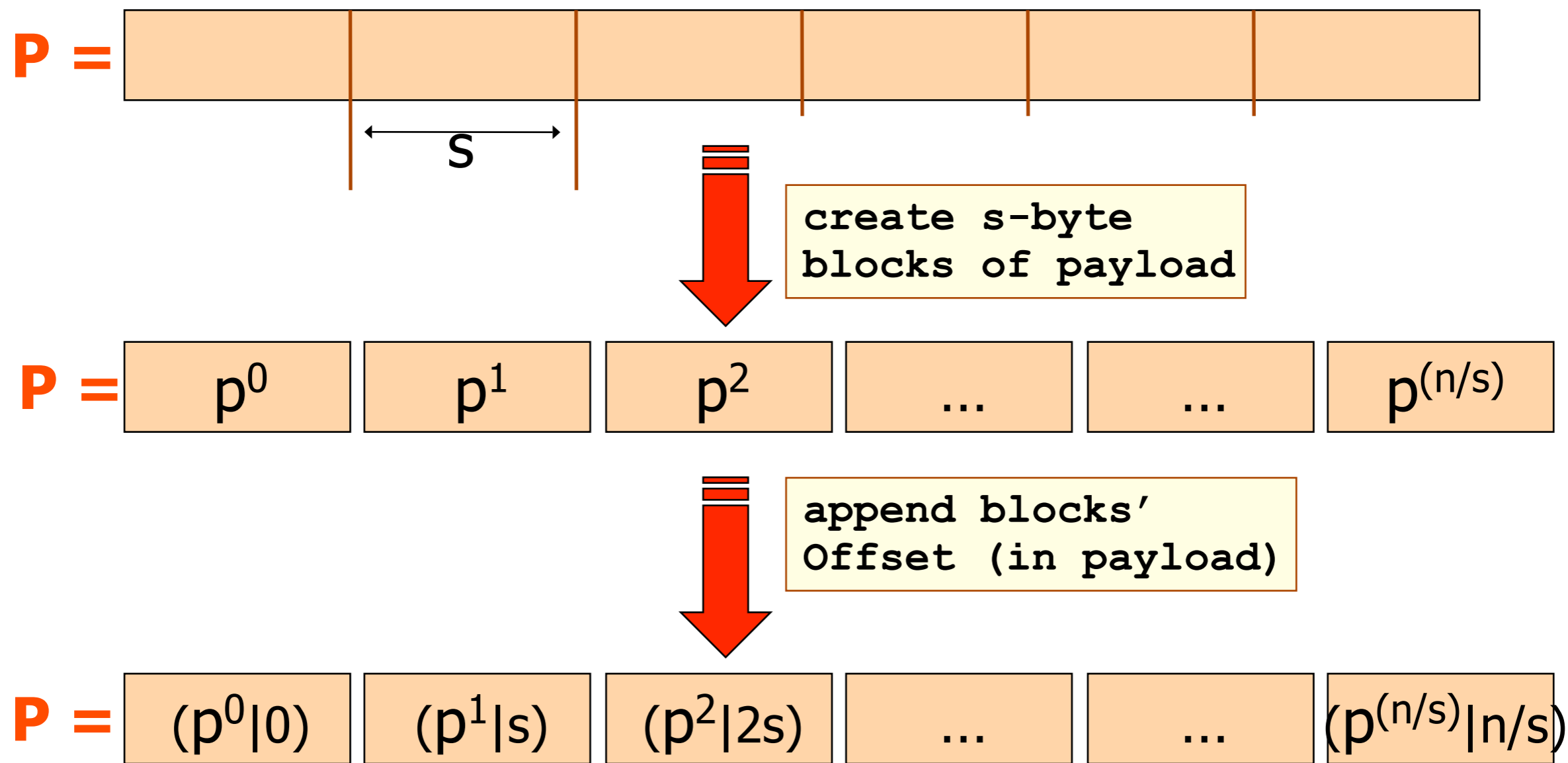  - n – number of elements in the set

# Packet Digests & Bloom Filters

- Space Efficient:

  - m/n=16 and 8 hashes (k=8)
    false positive (FP) = 5.74 x 10-4

  - No false negatives!

- However, suppose we don't have packets.

  - We only have some excerpts of payload

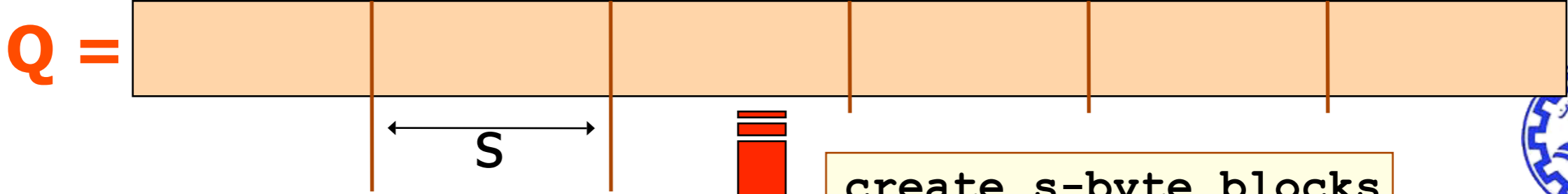  - Don't know where the excerpt was aligned in the packet

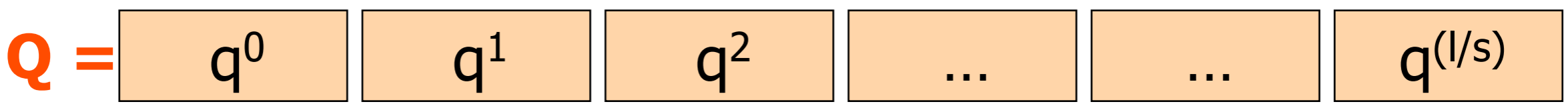**Extend Bloom Filters to support excerpt/substring matching**
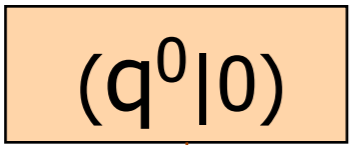
# Block-based Bloom Filter

$P =$ ⬜⬜⬜⬜⬜⬜

← $s$ →

**create s-byte blocks of payload**

$P =$ | $p^0$ | $p^1$ | $p^2$ | ... | ... | $p^{(n/s)}$ |

**append blocks' Offset (in payload)**

$P =$ | $(p^0|0)$ | $(p^1|s)$ | $(p^2|2s)$ | ... | ... | $(p^{(n/s)}|n/s)$ |

**Insert each block into a Bloom Filter**

$Q =$

create s-byte blocks
of query string

$Q =$ $q^0$ $q^1$ $q^2$ ... ... $q^{(l/s)}$

Try all possible offsets

$(q^0|0)$

$H_1(q^0|0)=1$
$H_2(q^0|0)=1$
$H_3(q^0|0)=0$

$(q^0|s)$ $(q^1|2s)$ $(q^2|3s)$

$q^0=p^1$ $q^1=p^2$ $q^2=p^3$

X

"$q^0q^1q^2$" was seen in
a payload at
offset 's'

[Memon]

**BBF =** $(p^0|0)$ $(p^1|s)$ $(p^2|2s)$ $(p^2|3s)$ ... $(p^{(n/s)}|n/s)$

**P1 =**

| A | B | R | A | C | A |
|---|---|---|---|---|---|

**P2 =**
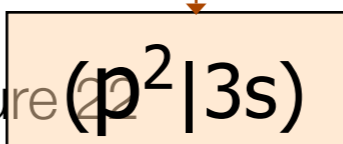
| C | D | A | B | R | A |
|---|---|---|---|---|---|

**BBF =**

| (A\|0) | (B\|s) | (R\|2s) | (A\|3s) | (C\|4s) | (A\|5s) |
|---|---|---|---|---|---|
| (C\|0) | (D\|s) | (A\|2s) | (B\|3s) | (R\|4s) | (A\|5s) |

## "Offset Collisions"

| (A\|0) | (B\|s) | (R\|2s) | (A\|3s) | (C\|4s) | (A\|5s) |
|---|---|---|---|---|---|
| (C\|0) | (D\|s) | (A\|2s) | (B\|3s) | (R\|4s) | (A\|5s) |

For query strings: "AD", "CB", "DR", "AA" etc. BBF falsely identifies them as seen in the payload!

Because BBF cannot distinguish between P1 and P2

[Memon]

# Hierarchical Bloom Filter

- **An HBF is basically a set of BBF for geometrically increasing sizes of blocks.**

**Hierarchical Bloom Filter**

**level-2** $(p^0 p^1 p^2 p^3 | 0)$

**level-1** $(p^0 p^1 | 0)$ $(p^2 p^3 | 2s)$ $(p^{(n/s-1)} p^{(n/s)} | (n/s-1))$

**level-0** $(p^0 | 0)$ $(p^1 | s)$ $(p^2 | 2s)$ ... ... $(p^{(n/s)} | n/s)$

**P =**

# Hierarchical Bloom Filter

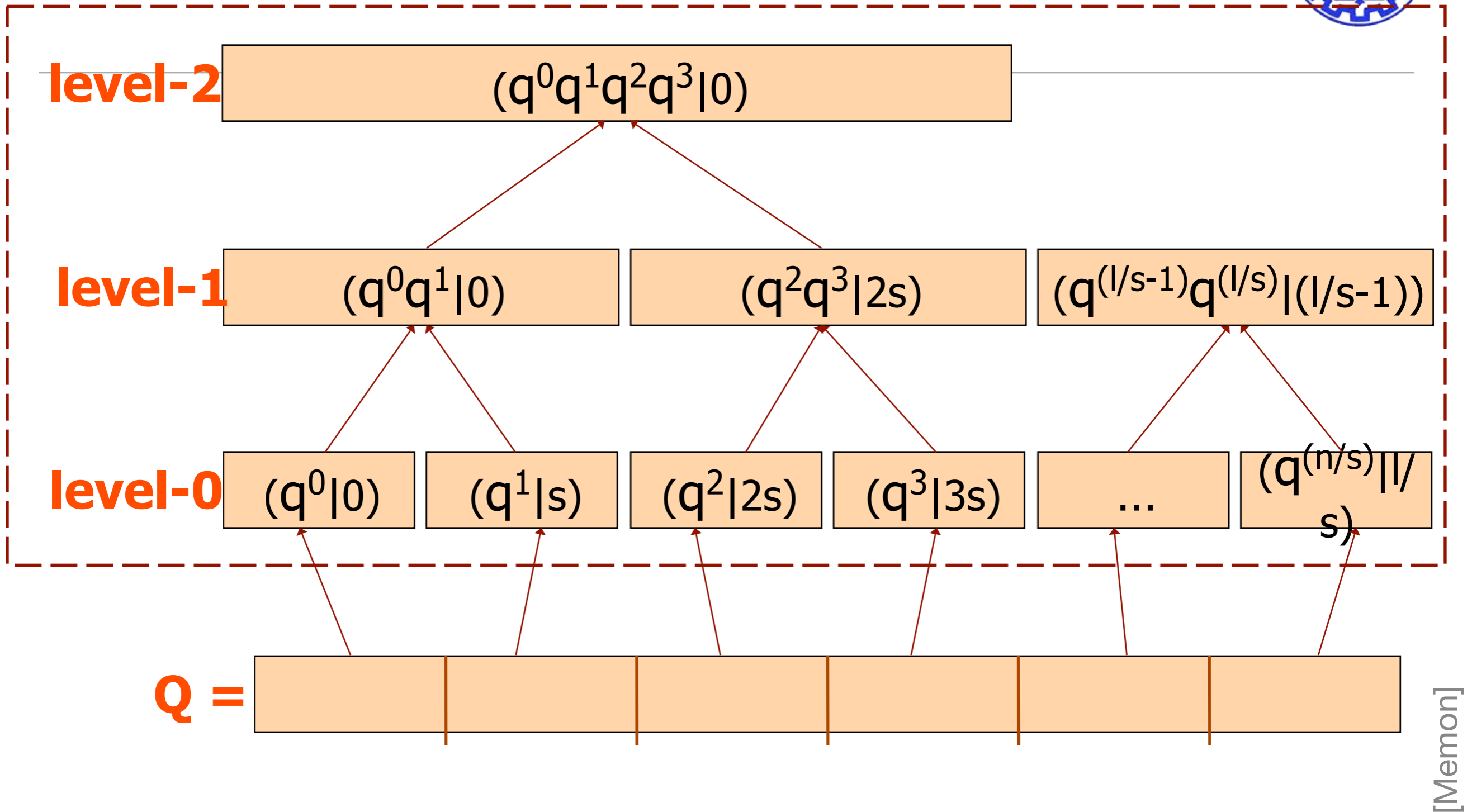**level-2** $(q^0 q^1 q^2 q^3 | 0)$

**level-1** $(q^0 q^1 | 0)$      $(q^2 q^3 | 2s)$      $(q^{(l/s-1)} q^{(l/s)} | (l/s-1))$

**level-0** $(q^0 | 0)$   $(q^1 | s)$   $(q^2 | 2s)$   $(q^3 | 3s)$   ...   $(q^{(n/s)} | l/s)$

**Q =**

- **Querying is similar to BBF.**
- **Matches at each level can be confirmed a level above.**
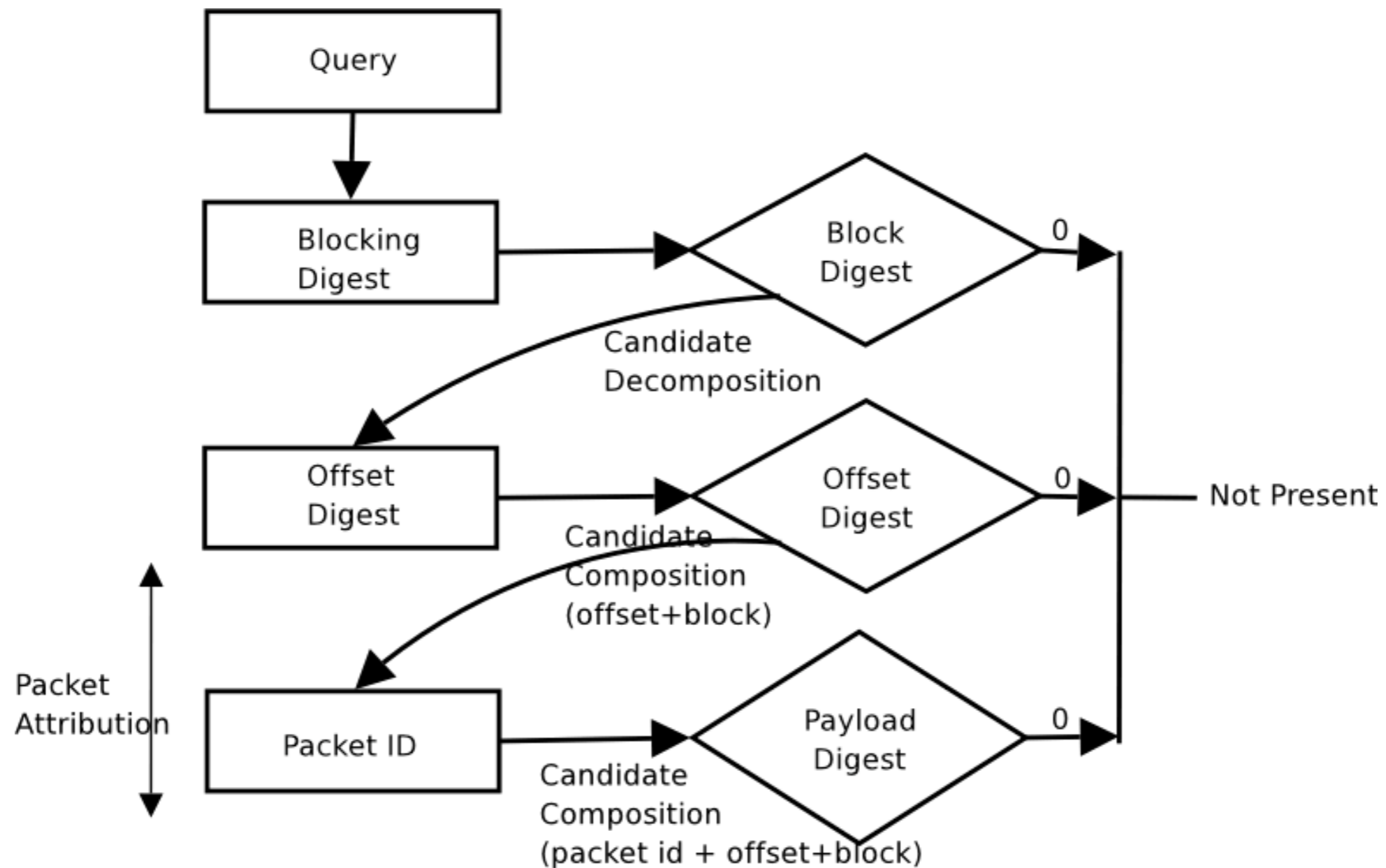
# Adapting an HBF for ForNet

- So far an HBF can attest for the presence of a bit-string in payloads

- We need to tie this bit-string to a source and/or destination hosts

- Our Approach:

  - Similar to tying an offset to a block/bit-string

  - In addition to inserting (block||offset) also insert (block||offset||hostid)
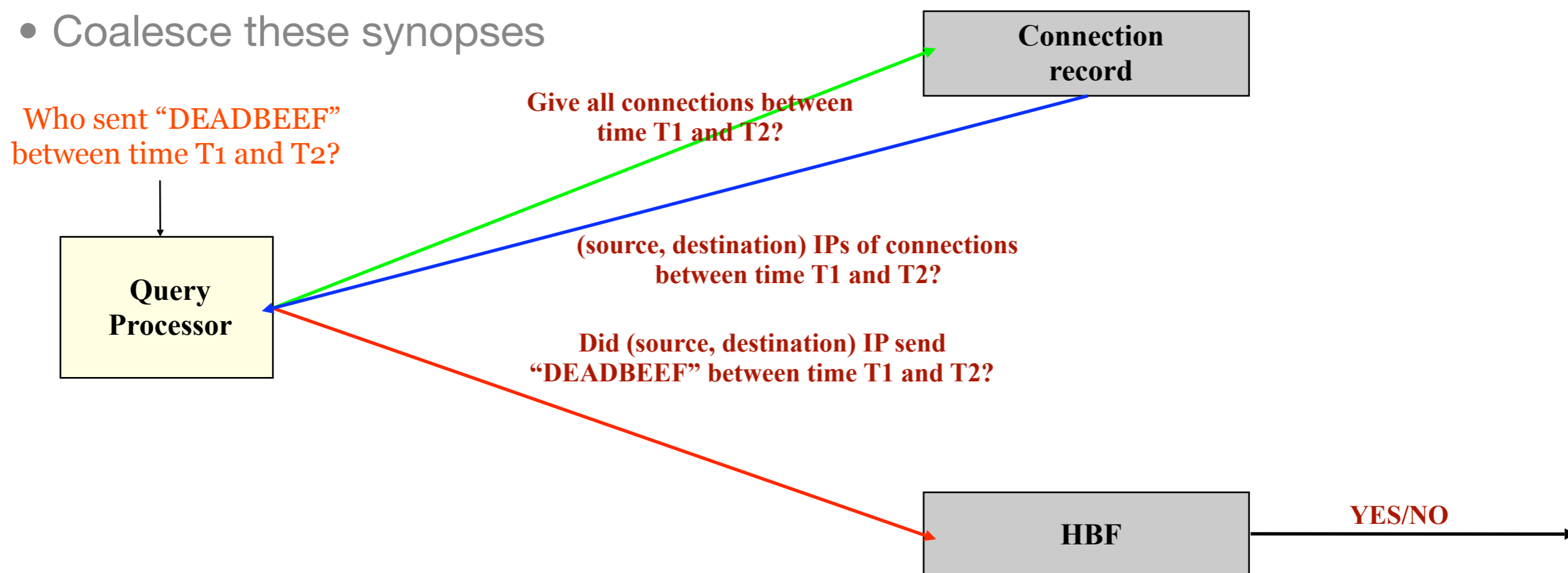
  - Hostid could be (srcIP||dstIP)

# How to run a query?

# Fusion of Synopses

- HBF requires:

  - Source IP, destination IP, excerpt

  - But where do we get Source IP, Destination IP

- Connection Record

  - Given two time intervals can give us list of source, destination IPs

- Coalesce these synopses

**Connection record**

Who sent "DEADBEEF" between time T1 and T2?

Give all connections between time T1 and T2?

**Query Processor**

(source, destination) IPs of connections between time T1 and T2?

Did (source, destination) IP send "DEADBEEF" between time T1 and T2?

**HBF**

YES/NO

# ForNet in Intranet Usage

- Determine victims of worm, trojans and other malware.

- Detection of potential victims of phising and spyware

# ForNet Deployed on Internet

- Traceback based on partial content of single packets

- Source of malware, worms, etc.

# Current Status

- Implemented a PC based SynApp device for placement within an intranet.

- Implemented Forensics Server with simple querying capabilities.

  - Current Forensic Server has 1.3TB of storage with over 3 months worth of data from the edge-router and two subnets

  - Normal bandwidth consumption of network is about a 1 – 2 TB/day

  - Synopses reduces this traffic to about 20GB/day
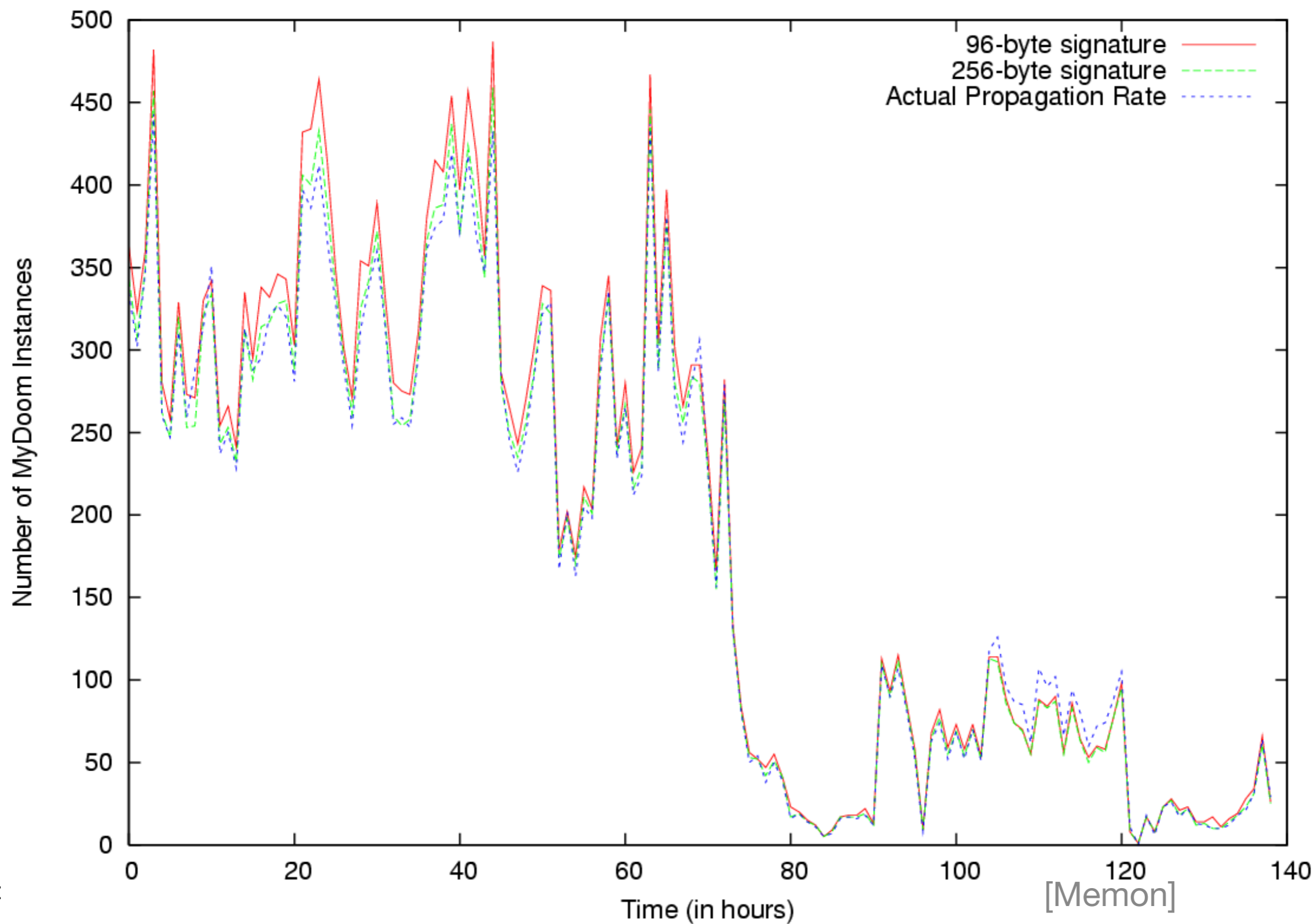
- Implemented Panorama (GUI Client)

# Tracking MyDoom

- Recorded all email traffic for a week
  - Using HBF and raw traffic
  - Was not aware of MyDoom during this collection
- When signatures became available we used them to query the system
  - To find hosts that are infected in our network
  - How the hosts were infected
- Some statistics:
  - 679 hosts originated at least one copy of the virus
    - 52 of which were in our network
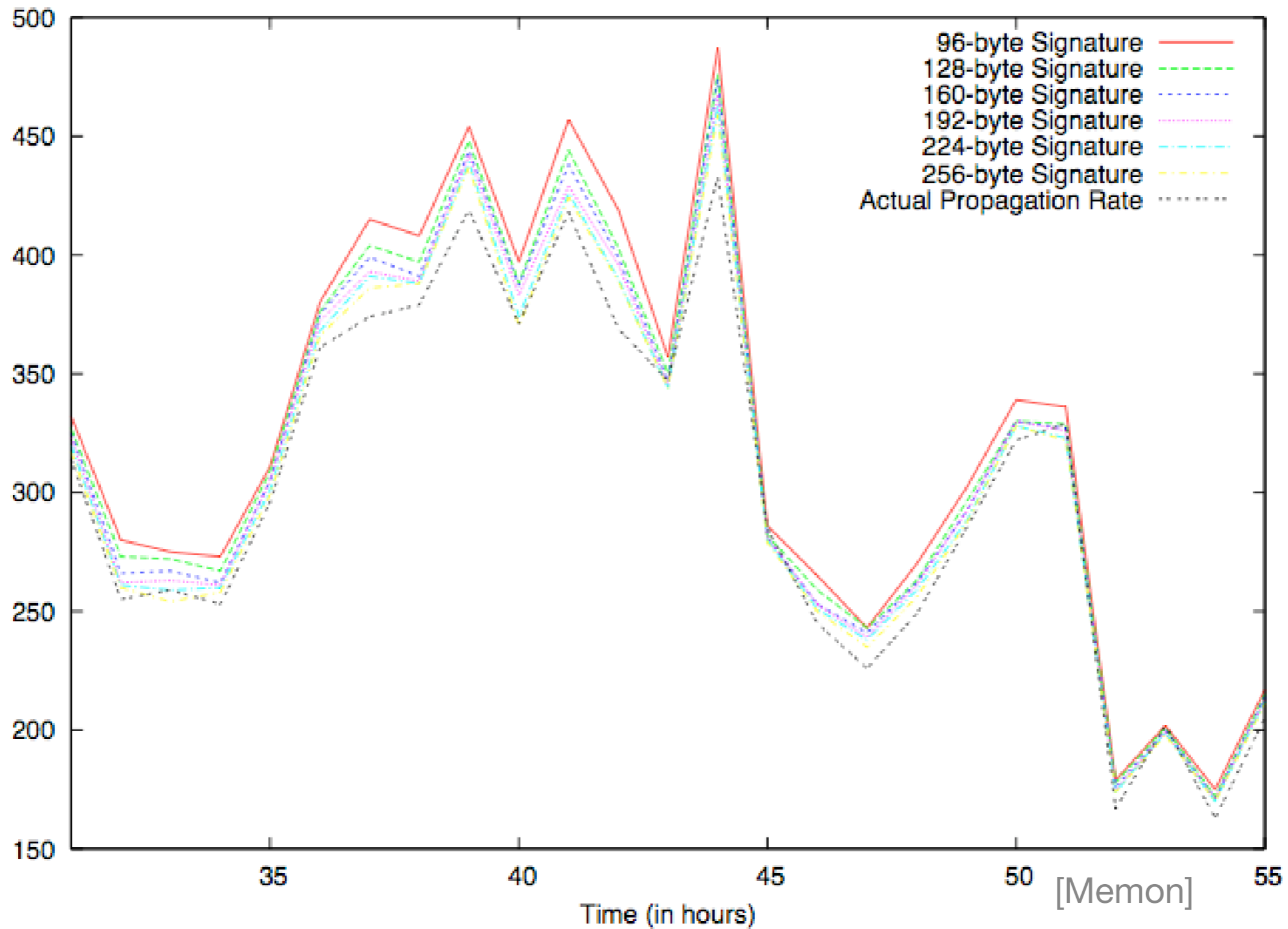  - These hosts sent out copies of the virus to 2011 hosts outside our network boundary

# Analyzing MyDoom Infections…

# MyDoom's One Day Zoom

# Acknowledgments/References

- [Xie] Toward a Framework for Internet Forensic Analysis, Yinglian Xie, presenting (Toward a Framework for Internet Forensic Analysis, V. Sekar, Y. Xie, D. Maltz, M. Reiter, H. Zhang, HotNets-III, 2004.) at the 100x100 clean slate prject presentatio in Pittsburgh, December 2004.

- [Memon] From slides prepared by Kulesh Shanmugasundaram, Hervé Brönnimann, and Nasir Memon, presented at various talks/lectures.