# Behavioral Detection of Stealthy Intruders

*Vern Paxson*

University of California, Santa Barbara

University of California, Berkeley

Georgia Institute of Technology

ARO/MURI Annual Review                                    September 9, 2011

# Overview of Recent UCB Efforts

- Empirical grounding of asset discovery & system roles/use/abuse in massive datasets
  - □ Drawn from operational environments, primarily Lawrence Berkeley National Laboratory (LBL)
    - 4K users, 12K hosts
- Scalable database technologies for archiving & querying against system event streams in real time
- *Behavioral-based detection of stealthy intruders*

# Finding Very Damaging Needles in Very Large Haystacks

- Motivation: some of the greatest threats to mission success arise from infiltrators unknown to have gained access to critical systems
- Particularly grievously damaging are long-term infiltrations that enable adversaries to develop a deep understand of mission components
- Such incidents might occur < 1/year …
  - … but cause more damage than all other intrusions combined

# Finding Very Damaging Needles con't

- Given event's very low frequency, adversary can expend extensive resources achieving initial compromise and maintaining a stealthy profile upon success
- Thus: behooves us to not focus on particular types of compromising attacks …
  - … but rather seek *behavioral indicators* that such an infiltration has occurred
    - Behavioral = look for signs of the presence of such an infiltrator
    - Can defend against very wide range of possible infiltration techniques including unknown ones ("zero days")

# Extracting Signals from Enormous Background Noise

- We can view this as a (highly nontraditional) signal processing problem:
  - □ Signal = behavioral indication of stealthy infiltrator
  - □ Domain = (incomplete) sensing extracted from site @ key locations
  - □ Noise = the huge amount of sensing that's <u>not</u> a stealthy intruder
- Key difference from traditional signal processing:
  - □ Data includes rich, highly <u>discontinuous</u> semantic structure
    - Both explicitly, and implicitly due to rules shaping activity (e.g., network protocols)
- We can leverage this semantic structure for much more powerful filtering than w/ generic approaches …
- … if we can determine the correct structural properties to exploit

# Seeking *Data At Scale* along with *Ground Truth*

- Crucial, underappreciated reality: conclusions derived from observing these semantic structures in simple environments (e.g., researcher's lab) do **NOT** "scale up" - they lack robustness
  - □ Reality of system use "in the wild" is much more complex & surprising than one expects
- In addition, we fundamentally require a degree of ground truth: the ability to determine the "right" answer for our haystack inferences
- LBL data provides the former due to scale …
- … and the latter due to our decades-long ties with its security & network operations staff

# Behaviors Associated With Stealthy Infiltrators

- Behavior #1: reconnaissance
  - Observation from analysis of past incidents: often upon subverting a system, attacker will investigate other systems reachable from it
  - Thus, associated *contact graph* should demonstrate "fan-out" or "depth" (contact = success *or failure* to log into a further host)
  - Approach: analyze site's interactive SSH traffic to infer multi-system access       (not yet working)
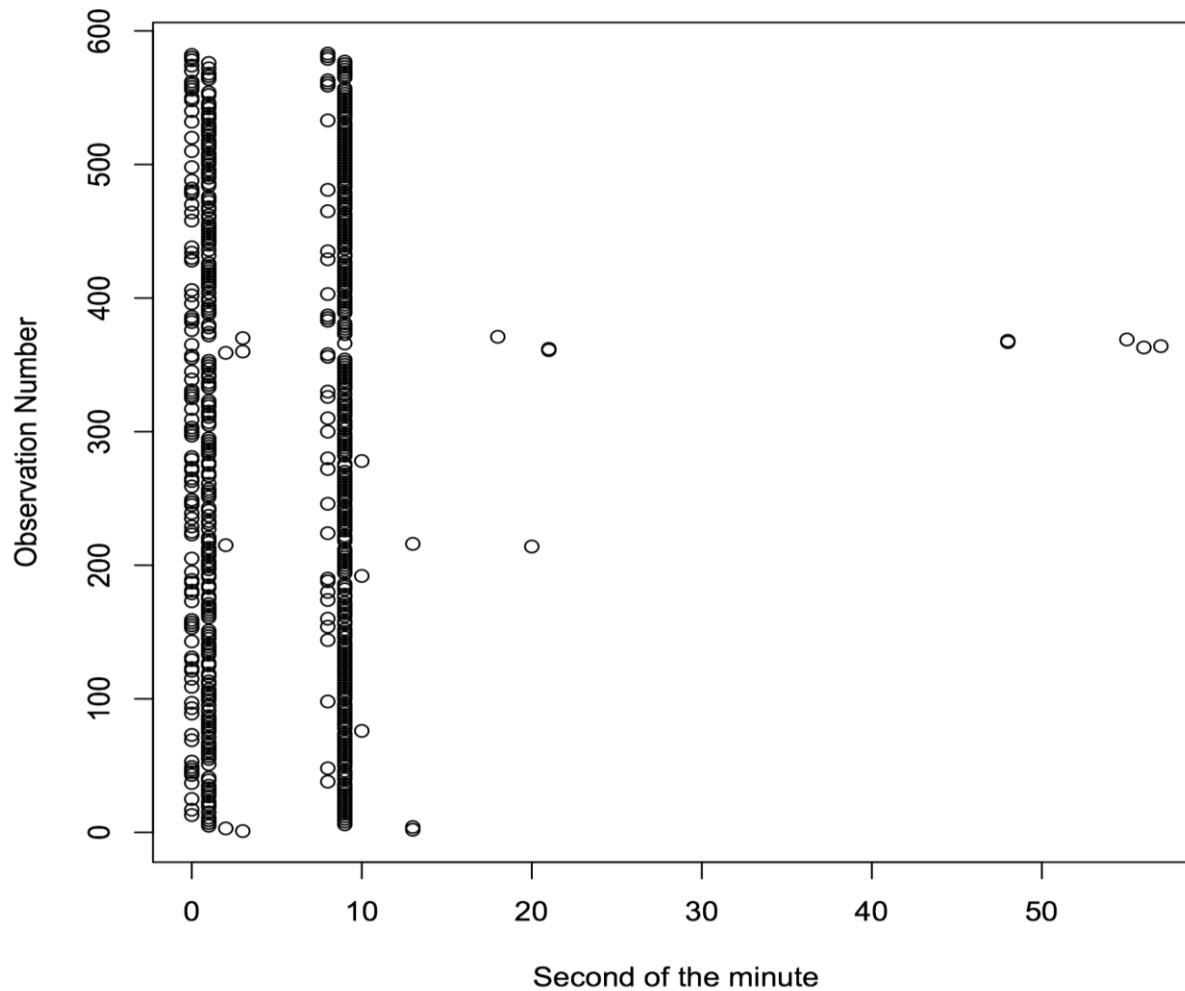- Behavior #2: covert tunneling
  - To access systems stealthily (and exfiltrate information unobserved), attackers can tunnel forbidden traffic inside another, benign/permitted protocol
  - Particularly attractive: DNS due to ubiquity
  - Approach: analyze DNS requests made by local systems that exhibits high entropy - many useful bytes transferred (working)

# Inferring Reconnaissance

- Initial data for analysis:
  - Logins made to LBL SSH servers from Jan 2009 - April 2011
    - Instrumented via *syslog*
      - Not comprehensive, but extensive (2K+ hosts, 3K+ accounts)
    - Data includes timestamp, originating host, server host, username, success
      - But <u>not</u> duration
    - 93.7 million records; most reflect internal logins
- Challenge #1: data is replete with automated activity
  - Not of interest for interactive reconnaissance
- Approach: sampling reveals that automation predominantly reflects periodic traffic …
  - … however NOT stationary
  - Insight: common periodicities align with per-minute/hour granularity

# Time of SSH Login for An LBL Client

QuickTime™ and a
decompressor
are needed to see this picture.

# $\chi^2$-based Testing for Automation

- Take series of activity timestamps, consider them mod 60 seconds or mod 60 minutes
- Place these in 6-60 bins (depending on amount of data available)
- Use $\chi^2$ to assess consistency with uniformity
- Failure $\Rightarrow$ automation candidate
- Now for remainder, compute size (depth/breadth) of potential contact graph
  - Quite problematic: lack of login durations, so how to tell if login from A to B overlaps with one from B to C?

# Focusing on Externally Initiated Traffic

- Key insight: most stealthy intruders begin their access with an *externally initiated* login …

- … and for those we do have connection durations
  - Due to monitoring by **Bro** system of site's border

- Provides sound upper bound on contact graph size

- 10-month assessment: 44 users had external logins w/ contact size > 2
  - Most were system administrators

- Asked operator to compare (w/o telling us!) against ground-truth database of known infiltrations
  - No matches :-(

- Now working on analyzing randomly selected incident

# Covert Channels for Communication & Exfiltration

- DNS lookups an integral part of Internet operation
- E.g., `www.cs.ucsb.edu` $\Rightarrow$ `128.111.41.37`
- What could be more benign?
- One of the few services ubiquitously allowed through firewalls
  - ☐ Though often restricted to site's designated local resolvers

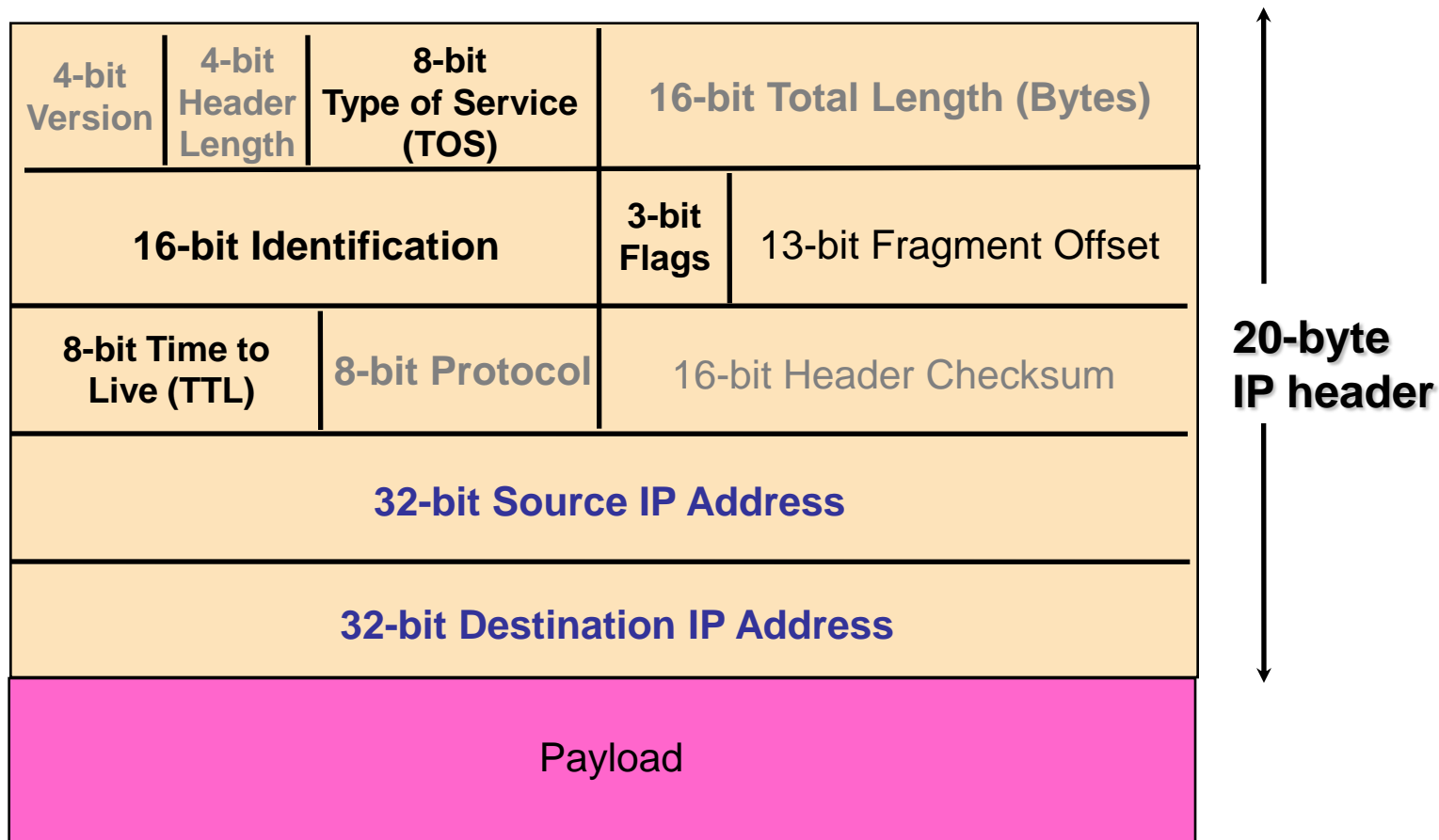- Exploiting this for arbitrary communication:

# Finding Room in DNS for Tunneling

"Questions" can include long names

As can Answers …

Even if query was for an A record (address), by returning a CNAME alias

| *16 bits* | *16 bits* |
|---|---|
| **Identification** | **Flags** |
| **# Questions** | **# Answer RRs** |
| **# Authority RRs** | **# Additional RRs** |
| **Questions** **(variable # of resource records)** | |
| **Answers** **(variable # of resource records)** | |
| **Authority** **(variable # of resource records)** | |
| **Additional information** **(variable # of resource records)** | |

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Payload | | | | |

**20-byte IP header**

*One way* to an IP packet inside a DNS packet:

```
version-4.hdrlen=5.TOS=0.len=81.<…etc…>.cs.ucsb.edu
```

Server can fully recover original IP packet, yet it's also a fully conformant DNS query

# How to Detect Such Tunneling?

- "Look for the funny name structure"
  - ☐ No good, attacker has <span style="color:red">enormous</span> degrees of freedom

- "Look for weirdly large lookups"
  - ☐ Here is where we encounter ~emergent behavior.  Huge traces have **benign** large lookups
  - ☐ E.g. 2.fnsroebjfvat-2seq-2sPuAao29aYJ1uoUqupzHgp 2uuqzSIRNRLzdDQVW6xNlb.ST9VNNN8IOEeFNNNO. fnsroebjfvat-pnpur.tbbtyr.pbz.u.00.s.sophosxl.net
    - This is actually antivirus software checking whether an executable is on a known blacklist

# Searching for DNS Tunnels in the Wild

- Data: 60 billion lookups from LBL (5.5 years)
  - Also developed scheme for working on border traffic
- Idea: Look for High-Entropy Domains
- E.g., "does `cs.ucsb.edu`" have a lot of different names looked up in it?
- Procedure:
  - Filter unique lookups per client
  - Identify and remove DNS search paths
    - Search path of foo.com generates zillions of XYZ.foo.com lookups
  - For every distinct domain, assess entropy of all lookups with it using gzip algorithm

# An Alternative Tunneling Approach: Repeated Lookups

- E.g., lookup `one.cs.ucsb.edu` to convey a 1 bit, `zero.cs.ucsb.edu` to convey a 0 bit
    - ☐ Obviously can use a larger codebook for efficiency
- Since we only consider unique lookups, `cs.ucsb.edu` would have entropy of ~7 bytes
- Insight: to detect, don't use unique lookups, let gzip algorithm do the work!
    - ☐ It already knows how to efficiently compress repeated instances of same symbol
- (Can do better still by ignoring repeats within TTL of last response)

# DNS Tunneling Detection Efficacy

- Threshold: inspect domains w/ entropy > 10KB
- Whitelist required
  - ☐ About 30 entries for LBL
  - ☐ E.g., `bl.barracuda.com`, `dnsbl.manitu.net`, `nerd.dk`
    - (ironically, most are themselves **blacklists**)

- With this in place, it works!
  - ☐ Found two such tunnels at LBL going back to May 2006
  - ☐ (Both represent intended use by staff)

# Summary of Accomplishments

- Very large-scale data regularization / calibration
  - \> 100 billion records of NetFlow / DHCP / LDAP / DNS / SSH / Email
- Developed behavioral detector for SSH-based reconnaissance
  - Including identification of automated SSH usage
- Developed behavioral detector for DNS-based tunneling
  - Successfully found previously unobserved tunnels
- VAST - Visibility Across Space & Time
  - Unified streaming database for tracking disparate forms of site activity
  - Limited progress last year due to student's very serious family situation
- *Towards Situational Awareness of Large-Scale Botnet Probing Events,* IEEE Transactions on Information Forensics & Security, 6(1), March 2011
  - Work mainly done prior to review year

# Plans for Subsequent Reviewing Periods

- Develop effective SSH-based behavioral detector by judiciously analyzing known LBL incidents
  - ☐ Extending then to per-user SSH usage models
- Refine DNS tunneling behavioral detector
  - ☐ Incorporate request timing
  - ☐ Develop models of asset interactions as inferred via DNS looups
- Mediate access to immense data store for refinement of asset inference algorithms
- Develop VAST streaming database to operational state
  - ☐ Deploy live in order to gather sensing data in real-time …
  - ☐ … and to support retrospective analysis / algorithm development

# Broader Issues

- **Collaboration plan**
  - Work with UCSB assessing asset/mission discovery protocols against LBL operational data
  - Incorporate this analysis into VAST's real-time streaming
- **Technology transition status**
  - Will port our successful detection algorithms (DNS tunneling now; SSH infiltrators hopefully soon) to the Bro network monitoring system (www.bro-ids.org)
    - Open-source NIDS/NIPS I developed
    - Runs at a dozen large sites: LBL, UCB, NCSA, OSU, GenDyn, …
    - Has major NSF infrastructure support grant
- **Personnel supported:**
  - Dr. Robin Sommer, staff scientist
  - Matthias Vallentin & Mobin Javed, UCB Ph.D. students