

# CAN WE IDENTIFY SPEAR PHISHING TARGETS BEFORE THE EMAIL IS SENT?

---

Jeremy Wendt, JD Doak,

Andy Wilson, Roger Suppona



**Sandia  
National  
Laboratories**



U.S. DEPARTMENT OF  
**ENERGY**



Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000..

# MOTIVATION: SPEAR PHISHING



*Attacker visits pages  
to find target data*



*Logs record  
all visit data*

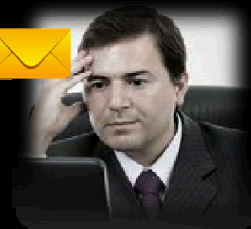
*Analysts identify and  
warn potential targets*



*Time Passes*

*Attacker crafts  
target-specific  
phishing email*

*Prepared user  
reports and deletes  
phishing email*



*Unprepared user  
follows link  
in email*

*Computer  
secure*



*Computer  
compromised*

# SOME IMPORTANT NOTES

- Spear Phishing moves at human speeds vs. network speeds
  - Attacker - identify targets, gather data, craft the email
    - hours to days
  - Victim - notice the email, read it, decide to click
    - seconds to hours/days
- Our Goal: Aid analysts in digesting logs

# SCALE

- Here at Sandia, we have a “medium sized” web-presence
    - 55 different machines serve web pages
    - 254 domain names served
    - ~2.5M distinct URLs found on SNL servers by crawling
    - 500K-1M entries per day
    - ~36K unique URLs requested per day
    - ~15K unique visitors per day
    - Downloading a single web page creates 1-20 entries in the log
      - HTML, images, CSS, JavaScript, etc.
  - And portions of the data are likely to be false
-

# DATA

- Each entry contains many pieces of data
  - Timestamp
  - Client IP
  - Client user agent string (UAS)
  - Requested URL
  - Refer string
  - X-Forwarded For (XFF)
  - Much more

# HELP THE ANALYST

- How can an analyst sift through such a mass of data quickly enough to find actionable data?
  - Improve the data, filter the data, sort the data, and present the data for better analyst triage
    - Distinguish crawler traffic from browser traffic
    - Sort the results so that most “interesting” sits at the top
    - Display the data so that a big picture is quickly visible

# DISTINGUISH CRAWLER FROM BROWSER

- Goal: Separate crawler traffic (e.g., Google indexer) from human-driven browser traffic
  - Some interesting other “patterns” could be hidden by bot “noise”
- UAS can be used to identify bot or not
  - But UAS is client-provided and can be falsified (or left empty)
  - Use UAS as initial grouping, look for other characteristics that distinguish the groups
    - Remember: Some of the clients are almost certainly lying
- NOTE: Visitors must leave 20 entries in the log to be included in this analysis
  - Really ~3-5 distinct webpages

# BOT CHARACTERISTICS

- Regular and fast: Get as many pages as quickly as you can!
- (Often) Polite: Don't hit a server too often or you'll kill it (and get blacklisted)
- Busy: Crawling the whole Internet takes a while
- (Often) Distributed: Only one machine needs to download each page
- Long memories: Don't redownload a resource until TIMEOUT passes





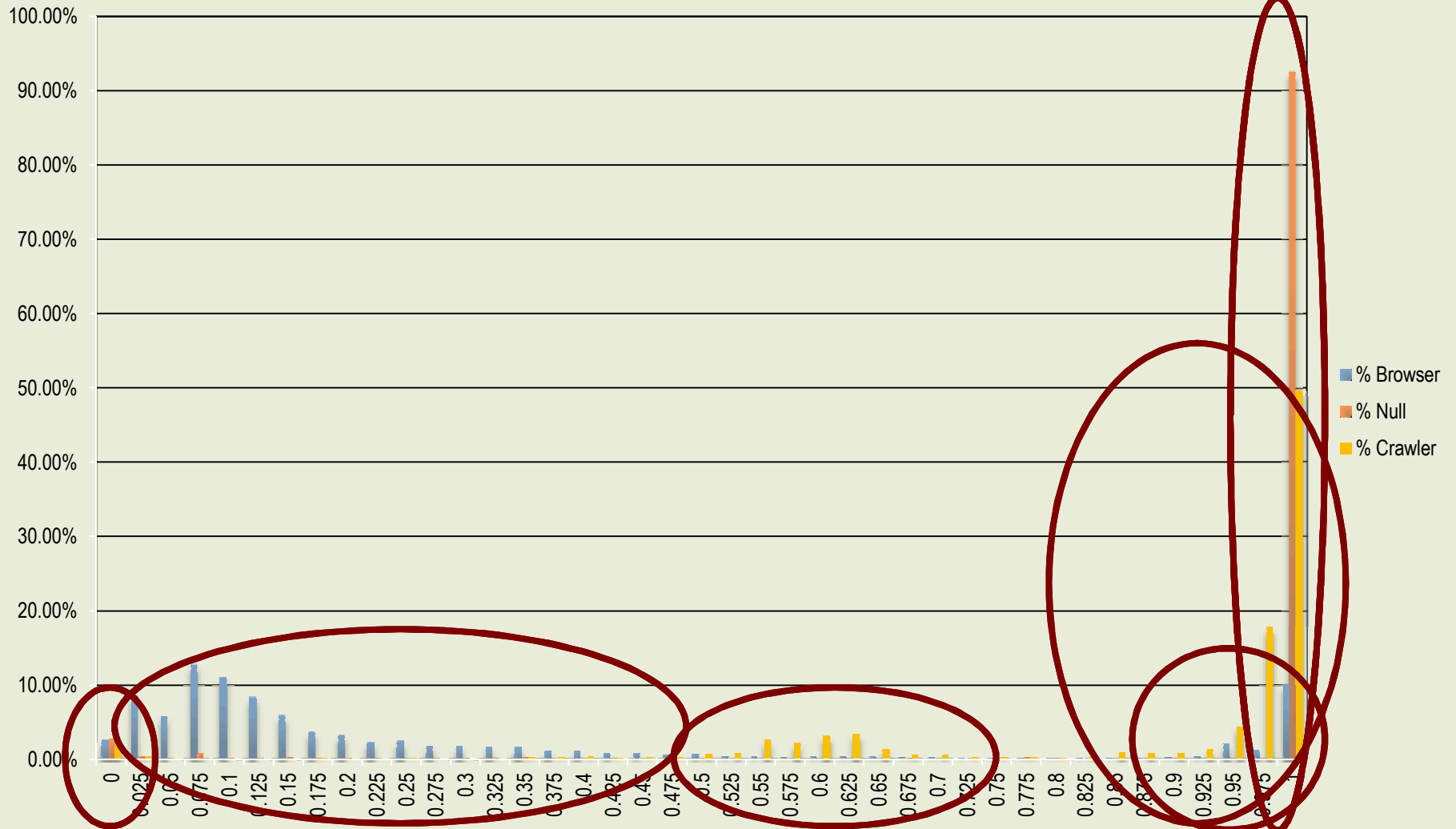
# BROWSER CHARACTERISTICS

- Bursty: Download all images/css/js/etc. for a page NOW!
  - Then do nothing until the user clicks again
- Lots of different file types
  - HTML/images/css/js/etc.
- Shorter memory: Will redownload content more regularly

# PERCENT HTML

- Browsers pull down 3-20 non-HTML documents for each HTML document to render the webpage
- Bots care most about indexing text (so often don't need images, etc.), and cache supporting documents longer than browsers
- $\%HTML = NumHTML / NumTotal$

# PERCENT HTML

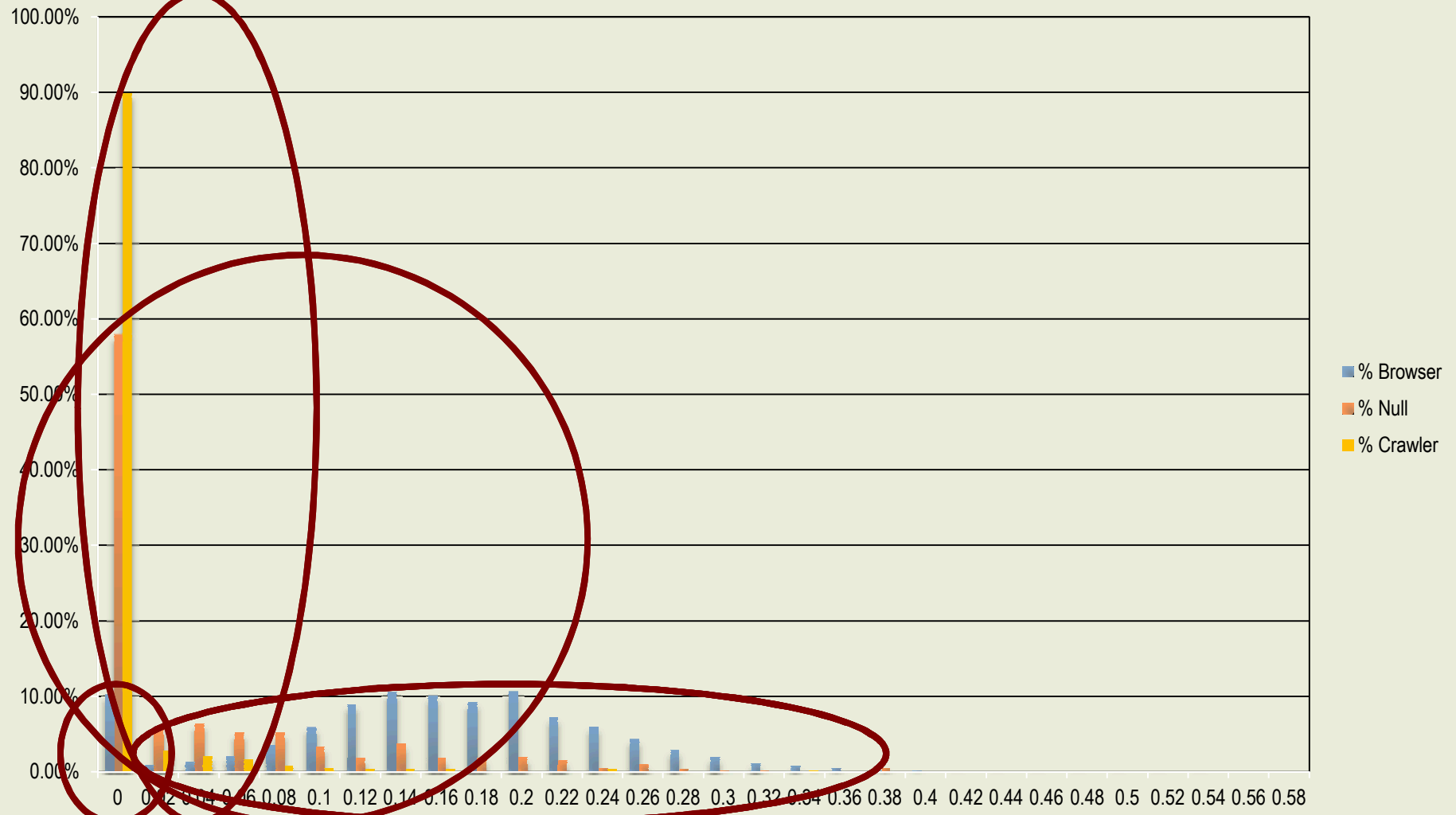


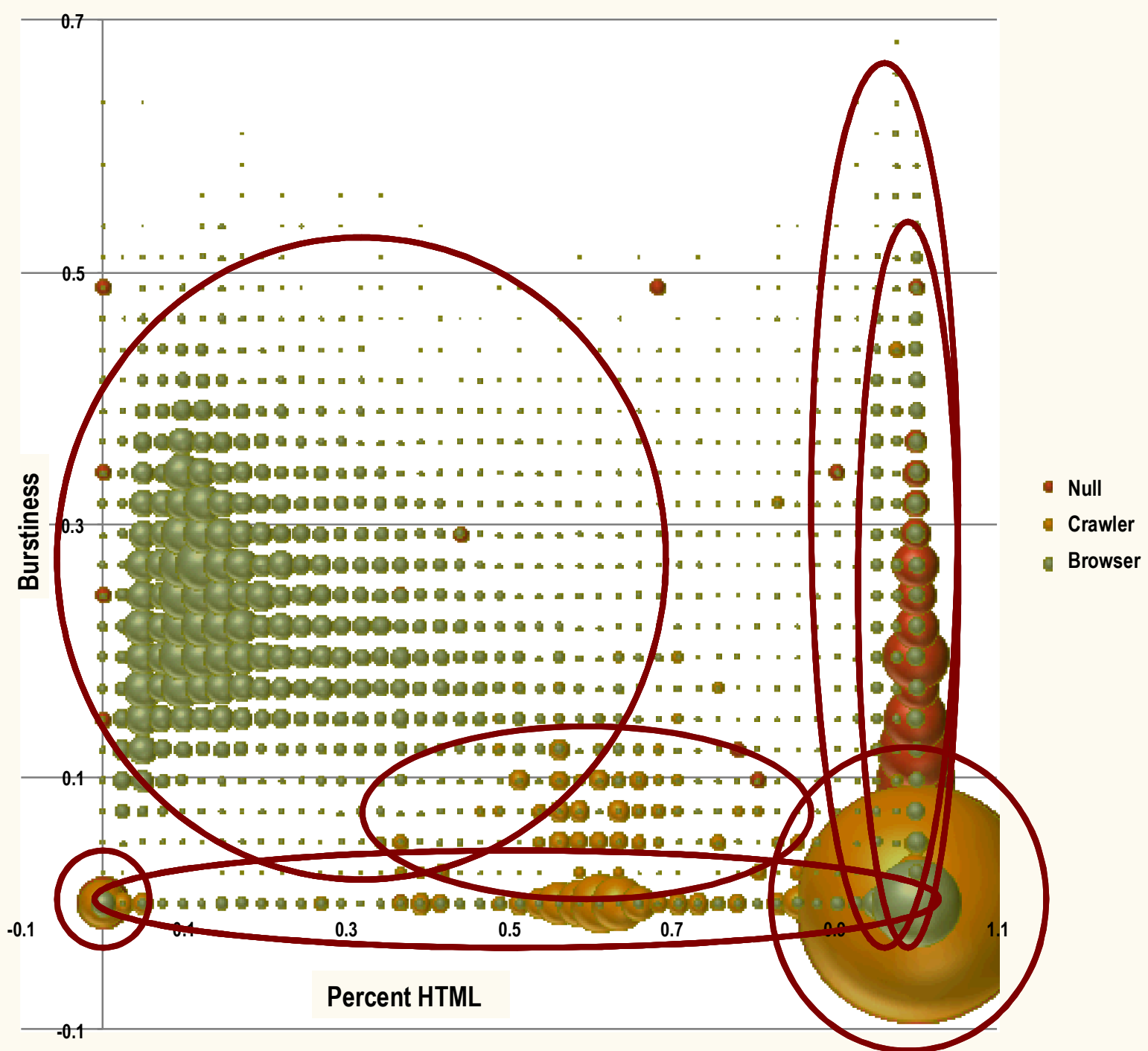
# BURSTS-TO-VISITS RATIO

- After browsers retrieve the HTML, they quickly parse it and request all supporting documents
- Bots only request pages from a site every second or so to keep from being blacklisted
- A “burst” is defined as more than N visits in an M second window by the same visitor
- $\text{Ratio} = \text{NumBursts} / \text{NumVisits}$



# BURST-TO-VISITS RATIO





# BOT-OR-NOT RESULTS

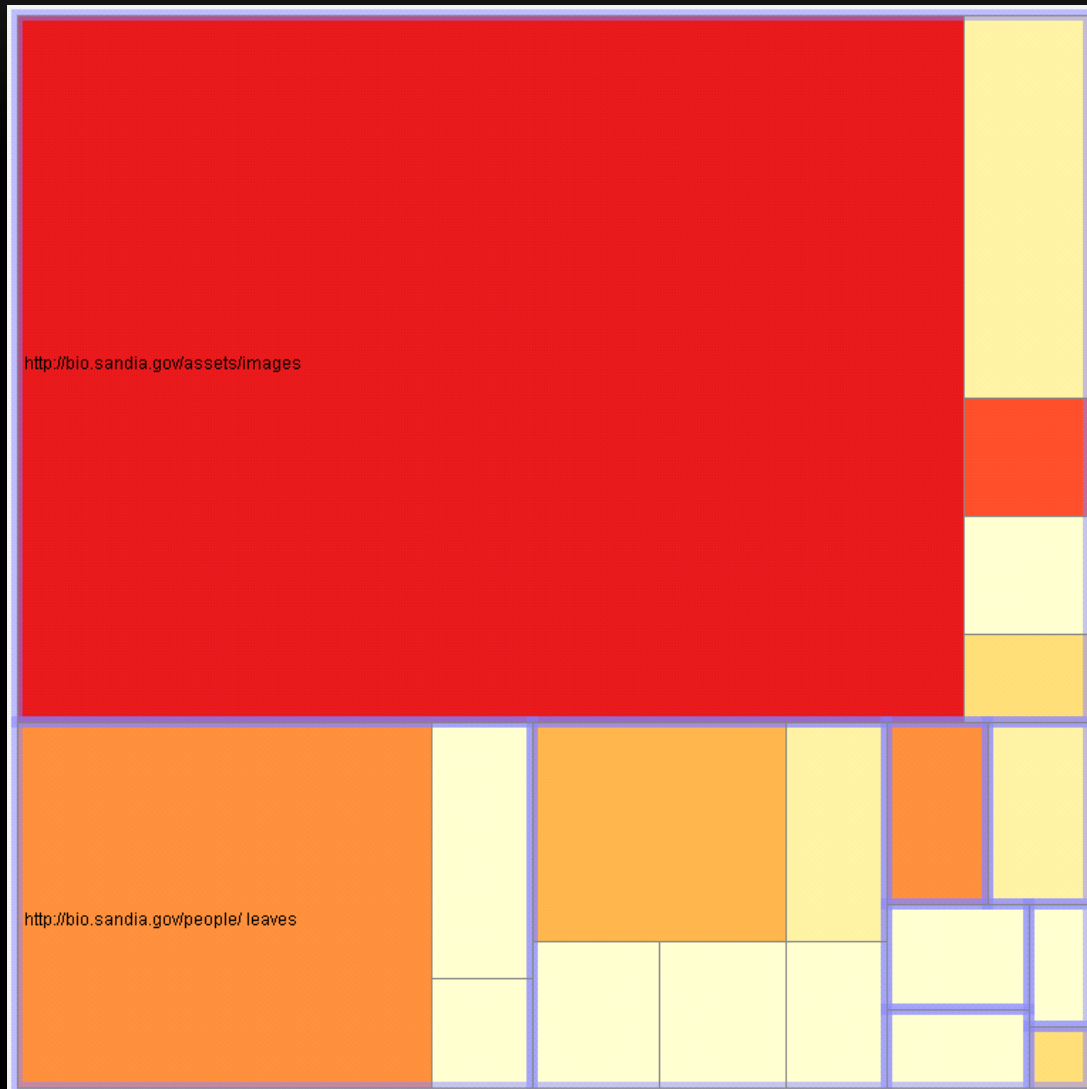
- Accuracy after combining the two features
  - UAS says “bot”: 98.6%
  - “null” UAS (assumed bot): 96.8%
    - NOTE: Evidence of ~1% browser-based null UAS
  - UAS says “browser”: 81.8%
    - NOTE: Evidence of 10-17% is bot

# BOT-OR-NOT RESULTS

- Visitors
    - 10.3K identified bots
    - 38.6K identified browsers
    - 166K too few visits
  - Visits
    - 3.2M made by bots
    - 3M made by browsers
    - 0.8M made by too few visits
  - Liars (possibly)
    - 8.7K say they're browser, but aren't
    - 30 (yes, thirty) say they're bot or null, but aren't
-



# TREEMAP



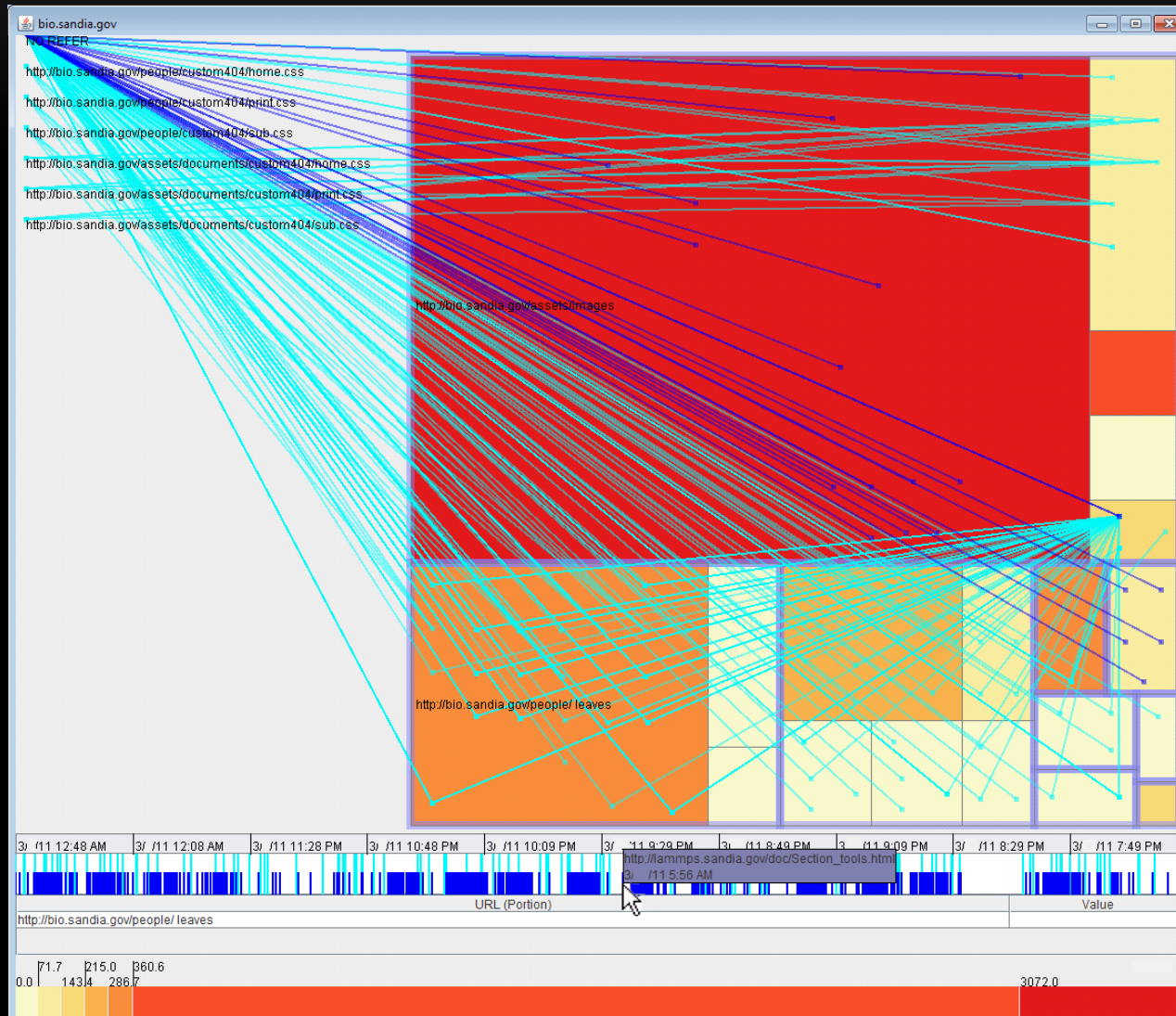
A rectangle's area represents the number of distinct webpages within each subdirectory.

A rectangle's color indicates the number of visits to webpages within each subdirectory.

The blue-outlined rectangles represent the level-one directories on [bio.sandia.gov](http://bio.sandia.gov).

The grey-outlined rectangles represent the level-two directories on [bio.sandia.gov](http://bio.sandia.gov).

# TRAVEL BY MAP



# QUESTIONS?



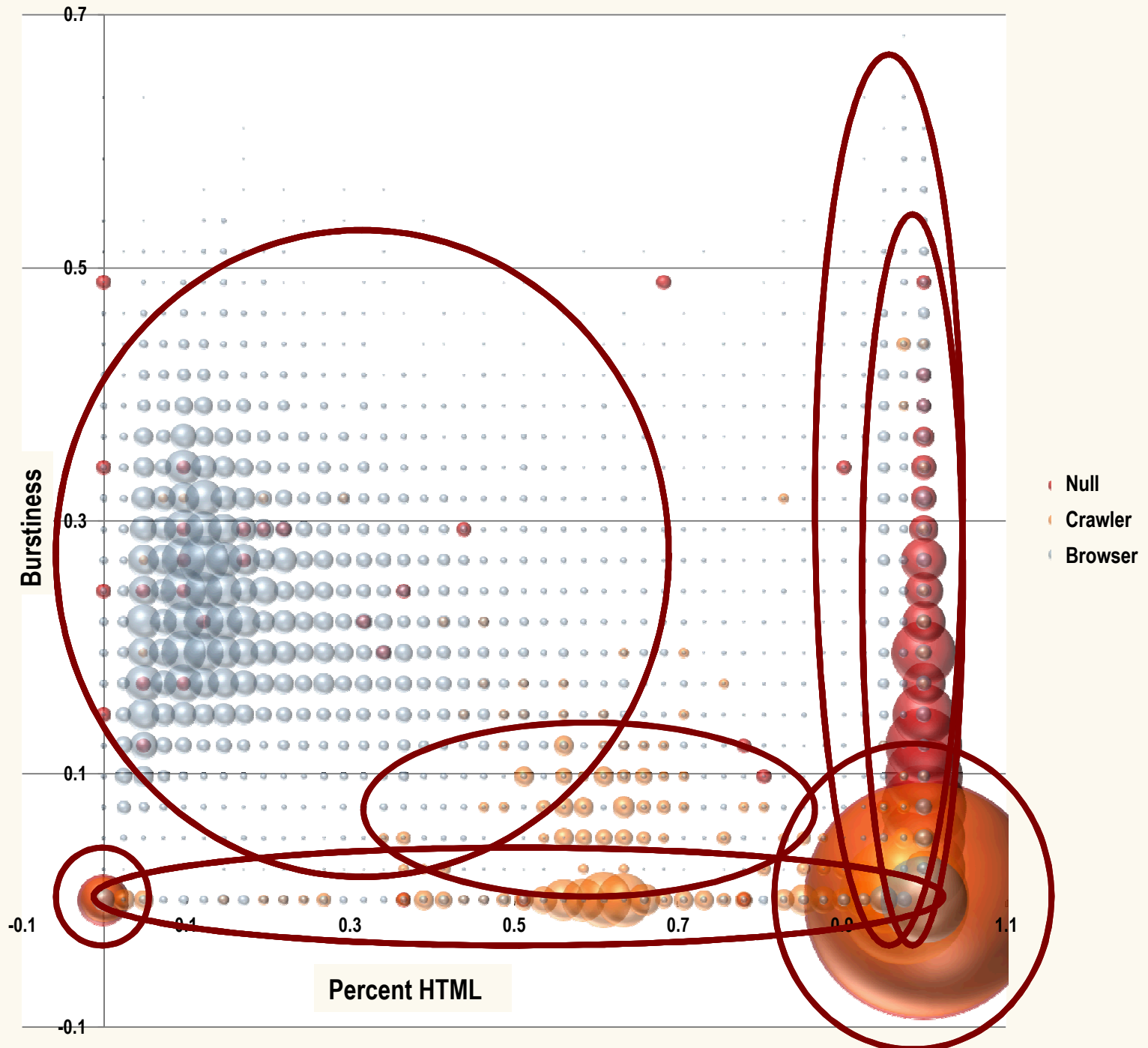
# BACKUP SLIDES



# CLIENT'S UNIQUE IDENTIFIER

- Due to DHCP, the same IP address can be used by different machines at different times
- Different applications on the same machine make HTTP requests for different reasons
  - UAS
- Proxied HTTP requests often use the last proxy server's IP address as the "client IP"
  - XFF (can contain private-space IPs which aren't globally unique)
- We decided to use the following pattern
  - "IP"/UAS
    - "IP" – If no XFF, use IP, else if 1<sup>st</sup> XFF addr is unique, use it, else use 1<sup>st</sup> XFF addr/1<sup>st</sup> unique address in XFF





# TRIAGE SUPPORT

The screenshot shows a web browser window with the title "Triage Example". The main content area displays a table with 10 columns. The first column contains IP addresses, and the subsequent columns contain various performance metrics. The table is sorted by the first metric column. The browser's address bar shows "http://energy.sandia.gov/wp-content...". The page content includes a list of links and a sidebar with a search bar and a list of categories.

IP Address	Browser/OS	0.005	0.008	0.991	0.991	0.182	0.667
198.40.41.250	Mozilla/4.0 (compatib...	0.005	0.008	0.991	0.991	0.182	0.667
69.198.131.153	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.667
69.127.72.30	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.667
98.207.224.185	Mozilla/5.0 (Macintos...	0.005	0.008	0.991	0.991	0.182	0.667
165.91.15.240	Mozilla/5.0 (Macintos...	0.001	0.001	0.992	0.992	0.286	0.4
72.102.78.227	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.4
50.53.212.133	Mozilla/4.0 (compatib...	0.005	0.008	0.991	0.991	0.182	0.5
71.170.172.203	Mozilla/4.0 (compatib...	0.003	0.006	0.989	0.989	0.4	0.182
74.105.112.147	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
209.19.33.193	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
68.40.207.169	Mozilla/5.0 (Macintos...	0.005	0.008	0.991	0.991	0.182	0.5
173.161.165.58	Mozilla/4.0 (compatib...	0.005	0.008	0.991	0.991	0.182	0.5
207.233.48.100	Mozilla/5.0 (Windows...	0.003	0.006	0.989	0.989	0.4	0.2
193.26.47.77	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
75.244.187.104	Mozilla/5.0 (Linux; U...	0.005	0.008	0.991	0.991	0.182	0.4
90.24.169.104	Mozilla/4.0 (compatib...	0.003	0.006	0.989	0.989	0.4	0.2
71.195.115.10	Mozilla/5.0 (compatib...	0.005	0.008	0.991	0.991	0.182	0.5
184.100.13.177	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
189.15.229.148	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
71.222.190.73	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.4
99.117.116.89	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
213.57.50.69	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
68.228.234.223	Mozilla/4.0 (compatib...	0.005	0.008	0.991	0.991	0.182	0.4
69.225.237.174	Mozilla/5.0 (Windows...	0.005	0.008	0.991	0.991	0.182	0.5
10.41.195.59/80.75...	Mozilla/4.0 (compatib...	0.005	0.008	0.991	0.991	0.182	0.5

The browser's address bar shows "http://energy.sandia.gov/wp-content...". The page content includes a list of links and a sidebar with a search bar and a list of categories.