

Beyond Blacklists: Learning to Detect Malicious Web Sites from Suspicious URLs

Justin Ma, Lawrence Saul, Stefan Savage, Geoff Voelker
Computer Science & Engineering
UC San Diego

Presentation for KDD 2009

June 30, 2009

Detecting Malicious Web Sites

URL = Uniform Resource Locator

<http://www.bfuduuioo1fp.mobi/ws/ebayisapi.dll>

<http://fblight.com>

<http://mail.ru>

<http://www.sigkdd.org/kdd2009/index.html>

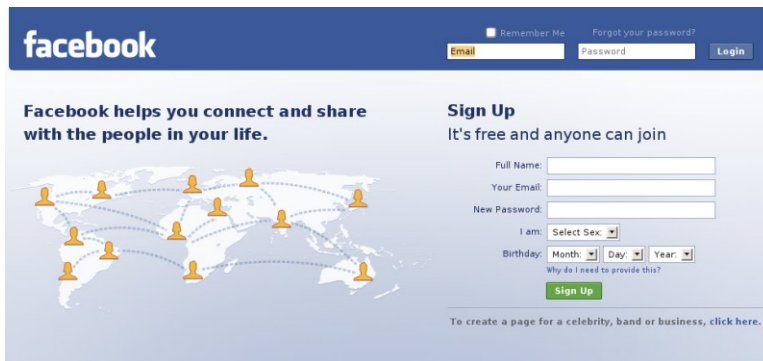
**safe without
committing to
risky actions**

The screenshot shows an email interface. At the top, it says "Wake UP your Feeling". The email is from "order@dmconcepts.com" and is dated "Mar 17 (10 days ago)". The main body of the email contains a red box with the text: "May be you need to try it" and "It works almost immediately". A red circle is drawn around this text. Below the red box, there are "Reply" and "Forward" buttons. At the bottom of the email, there is a footer with contact information: "You are receiving this newsletter because you subscribed to the About Today Newsletter at http://www.about.com/nl/usgs.htm?nl=todays&e=jtma@cs.ucsd.edu", "About respects your privacy. Our Privacy Policy.", "Our Contact Information.", "249 West 17th Street", "New York, NY, 10011".

Problem in a Nutshell

- URL features to identify malicious Web sites
 - No context, no content
- Different classes of URLs
 - Benign, spam, phishing, exploits, scams...
 - For now, distinguish benign vs. malicious

facebook.com



fblight.com



State of the Practice

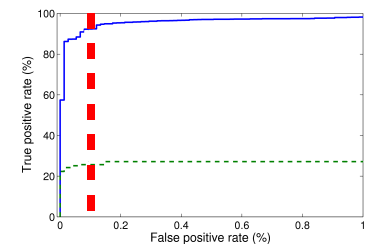
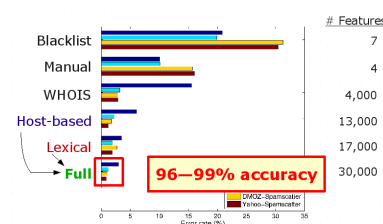
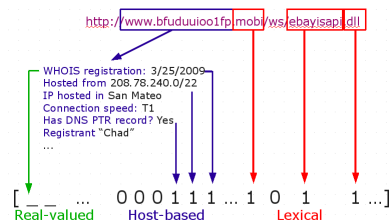
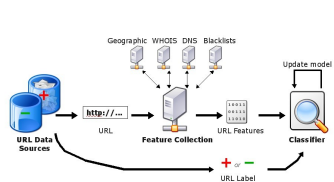
- Current approaches
 - Blacklists [SORBS, URIBL, SURBL, Spamhaus]
 - Learning on hand-tuned features [Garera et al, 2007]
- Limitations
 - Cannot predict **unlisted sites**
 - Cannot account for **new features**
- Arms race



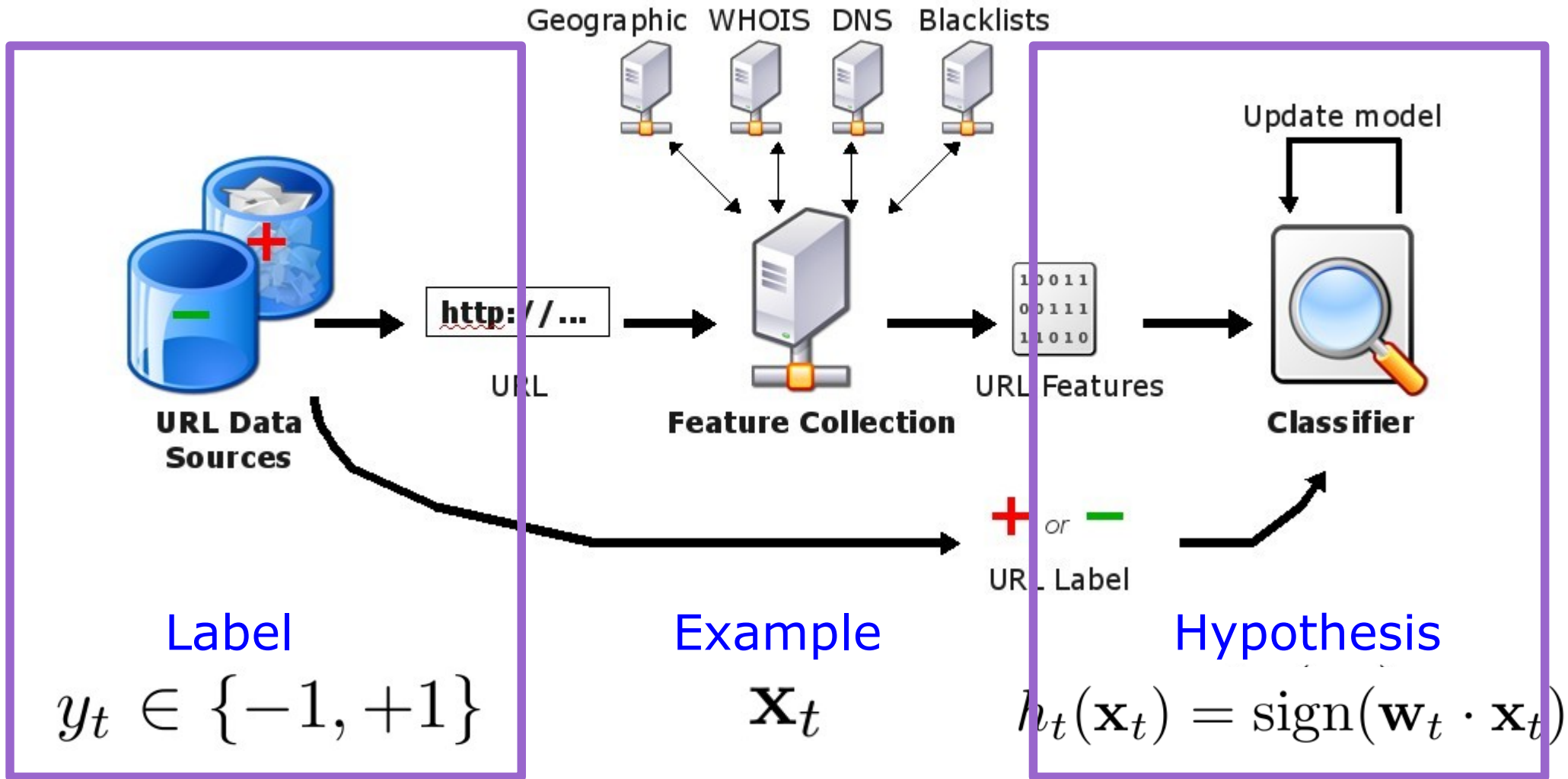
More automated approach?

Today's Talk

- Motivation
- System overview
 - Training data
 - Algorithms
 - **Features** ← focus of today's talk
- Experimental results
- Conclusion



URL Classification System



Data Sets

- Malicious URLs
 - 5,000 from PhishTank (phishing)
 - 15,000 from Spamscatter (spam, phishing, etc)
- Benign URLs
 - 15,000 from Yahoo Web directory
 - 15,000 from DMOZ directory
- Malicious x Benign → 4 Data Sets
 - 30,000 – 55,000 features per data set

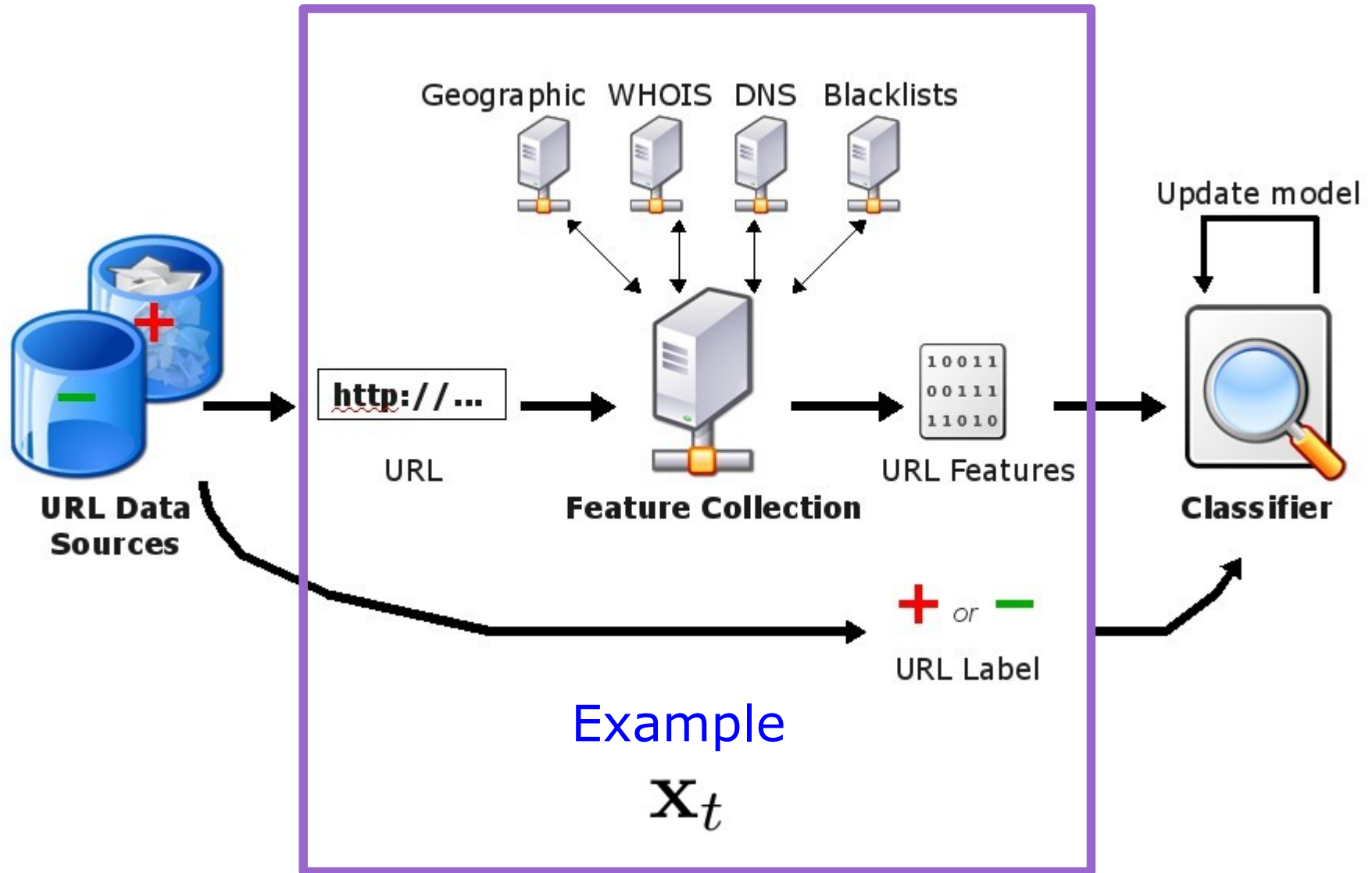
Algorithms

- Logistic regression w/ L1-norm regularization

$$L(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \log P(y_i | \mathbf{x}_i, \mathbf{w}) - \lambda \|\mathbf{w}\|_1$$

- Implicit feature selection
 - Easier to interpret
-
- Other models
 - Naive Bayes
 - Support vector machines (linear, RBF kernels)

Today's Focus



Feature vector construction

http://www.bfuduuiioo1fp mobi/ws/ebayisapi dll

WHOIS registration: 3/25/2009
Hosted from 208.78.240.0/22
IP hosted in San Mateo
Connection speed: T1
Has DNS PTR record? Yes
Registrant "Chad"
...

[_ _ ... 0 0 0 1 1 1 ... 1 0 1 1 ...]

Real-valued Host-based Lexical

Features to consider?

1) Blacklists



2) Simple heuristics



3) Domain name registration



4) Host properties



5) Lexical



(1) Blacklist Queries

- List of known malicious sites
- Providers: SORBS, URIBL, SURBL, Spamhaus

Blacklist queries as features

<http://www.bfuduuioo1fp.mobi>

In blacklist?

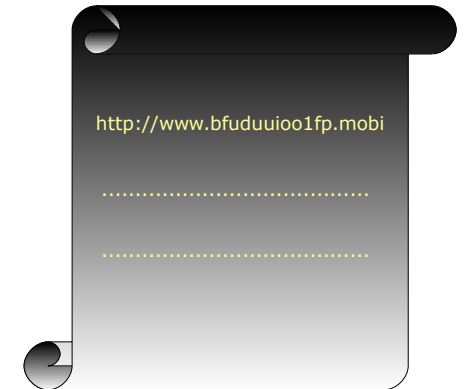


Yes

<http://fblight.com>



No



(2) Manually-Selected Features

[Fette et al., 2007][Zhang et al., 2007][Bergholz et al., 2008]

- Considered by previous studies
 - IP address in hostname?
 - Number of dots in URL
 - WHOIS (domain name) registration date

<http://72.23.5.122/www.bankofamerica.com/>

<http://www.bankofamerica.com.qytrpbcw.stopgap.cn/>



stopgap.cn registered
28 June 2009

(3) WHOIS Features

- Domain name registration
 - Date of registration, update, expiration
 - Registrant: Who registered domain?
 - Registrar: Who manages registration?

<http://yammeringyellowtail.com>

<http://angryalbacore.com>

<http://sleazysalmon.com>

<http://mangymackerel.com>

Registered on
29 June 2009
By **SpamMedia**



(4) Host-Based Features

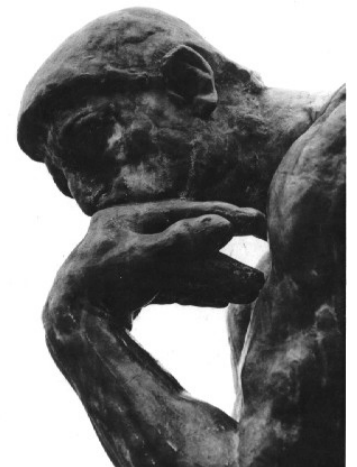
- Blacklisted? (SORBS, URIBL, SURBL, Spamhaus)
- WHOIS: registrar, registrant, dates
- IP address: Which ASes/IP prefixes?
- DNS: TTL? PTR record exists/resolves?
- Geography-related: Locale? Connection speed?



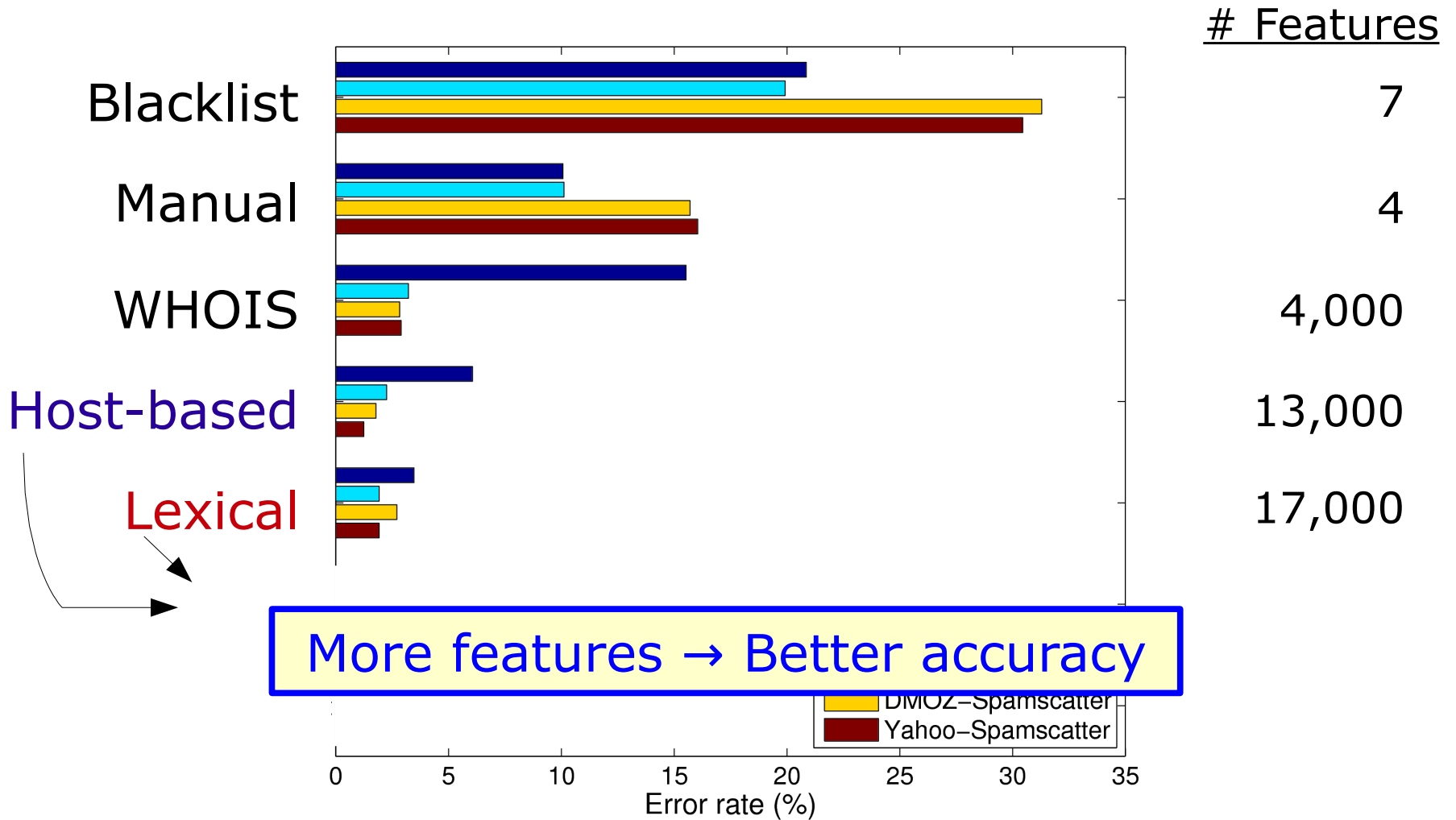
(5) Lexical Features

- Tokens in URL hostname + path
- Length of URL
- Number of dots

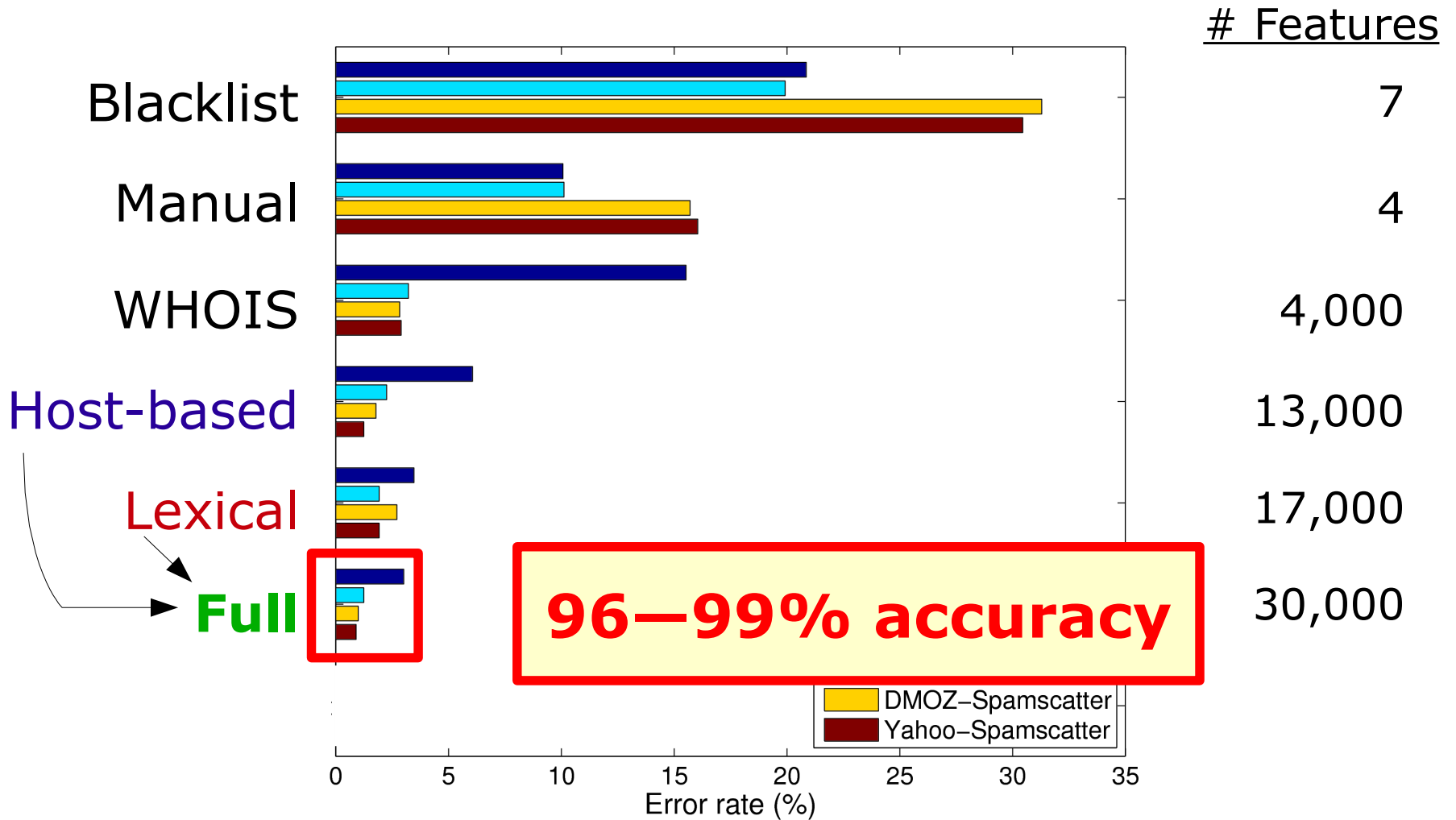
http://www.bfuduuioo1fp.mobi/ws/ebayisapi.dll



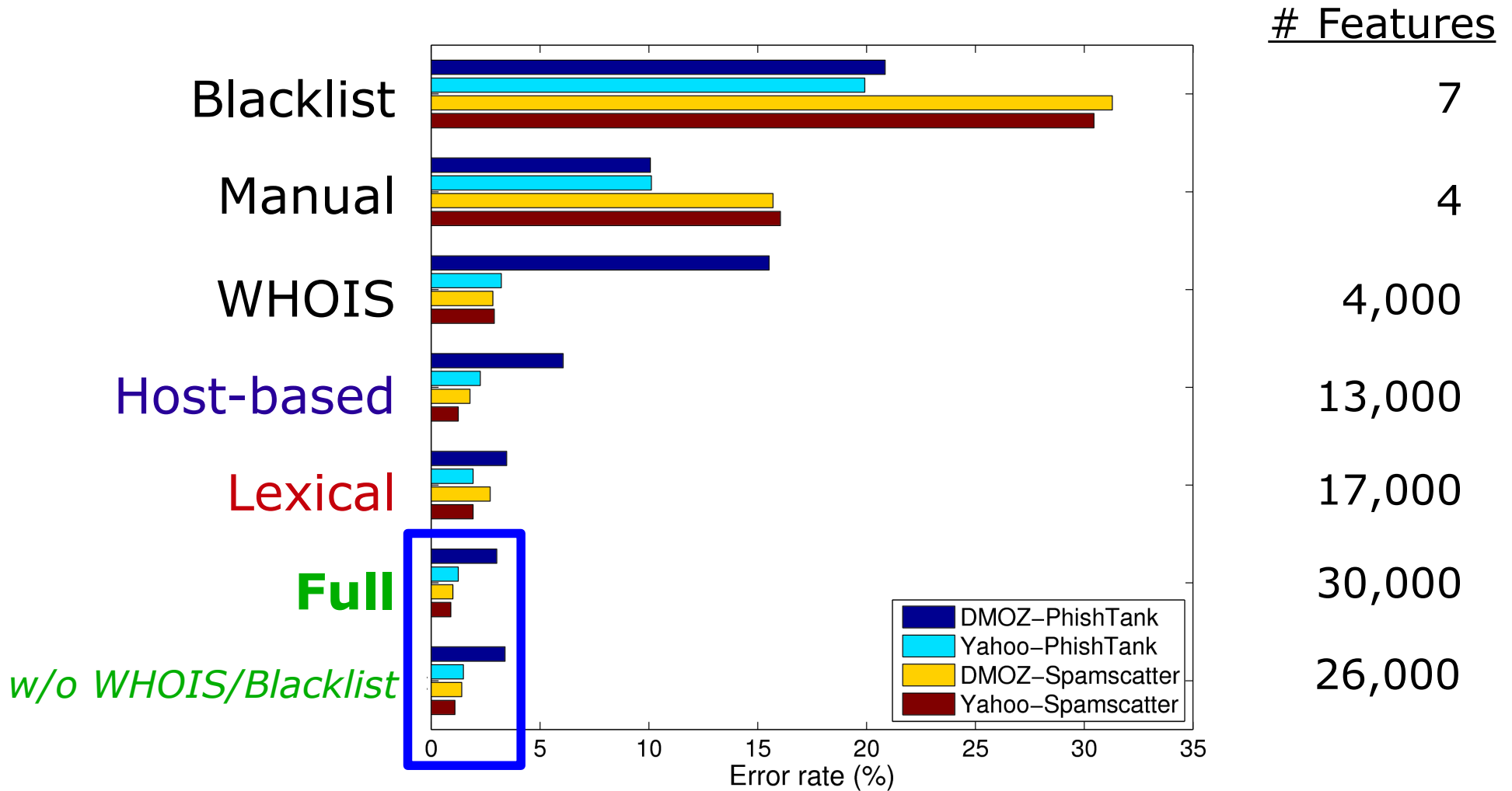
Which feature sets?



Which feature sets?

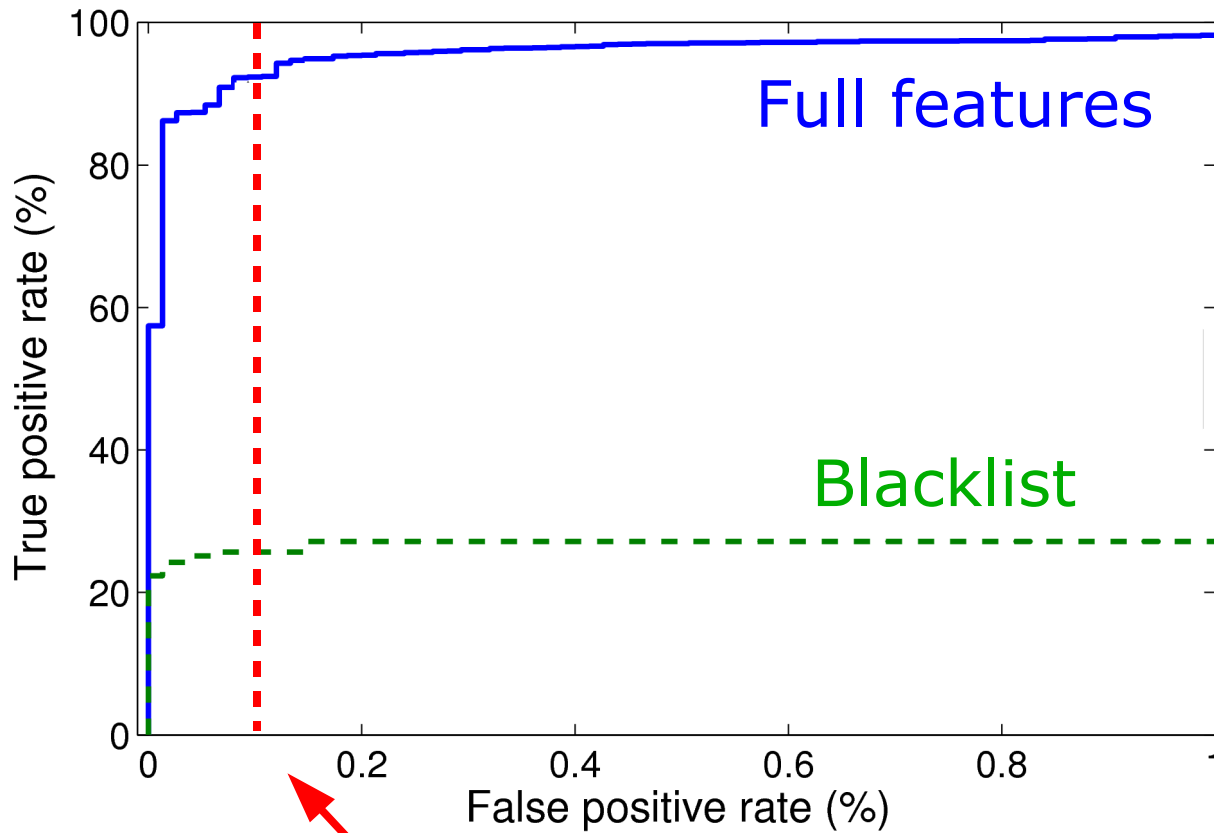


Which feature sets?



Beyond Blacklists

Yahoo-PhishTank



Higher detection rate for given false positive rate

Limitations

- **False positives**
 - Sites hosted in disreputable ISP
 - Guilt by association
- **False negatives**
 - Compromised sites
 - Free hosting sites
 - Redirection (but we consider TinyURL malicious :)
 - Hosted in reputable ISP
- **Future work:** Web page content

Conclusion

- Detect malicious URLs with high accuracy
 - Only using URL
 - Diverse feature set helps: 99% w/ 30,000+ features
 - Model analysis (more in paper)
- Our related efforts
 - Online learning for URL reputation [ICML 2009]
- Future work
 - Scaling up for deployment