

# Harnessing computing power for security: BitBlaze, WebBlaze, and real-time spam URL filtering

Devdatta Akhawe, Domagoj Babić, Adam Barth, Juan Caballero,  
Chris Grier, Steve Hanna, Justin Ma, Lorenzo Martignoni,  
Stephen McCamant, Feng Mao, James Newsome, Vern Paxson,  
Prateek Saxena, Dawn Song, and Kurt Thomas  
{smcc,dawnsong}@cs.berkeley.edu

University of California, Berkeley

# Computer security: bad news

## Apple slammed over iPhone, iPad location tracking

By JORDAN ROBERTSON AP Technology Writer

## Online scammers jump on bin Laden news

By BARBARA ORTUTAY AP Technology Writer

## Data breach at security firm linked to attack on Lockheed

By Christopher Drew and John Markoff

New York Times

Posted: 04/26/2011 08:50:24 PM PDT

## Sony: Credit data risked in PlayStation outage

By RYAN NAKASHIMA and JORDAN ROBERTSON AP Business Writers

Posted: 04/26/2011 02:15:27 PM PDT

Updated: 04/26/2011 08:50:24 PM PDT

More powerful computers seem to just increase our exposure to security threats

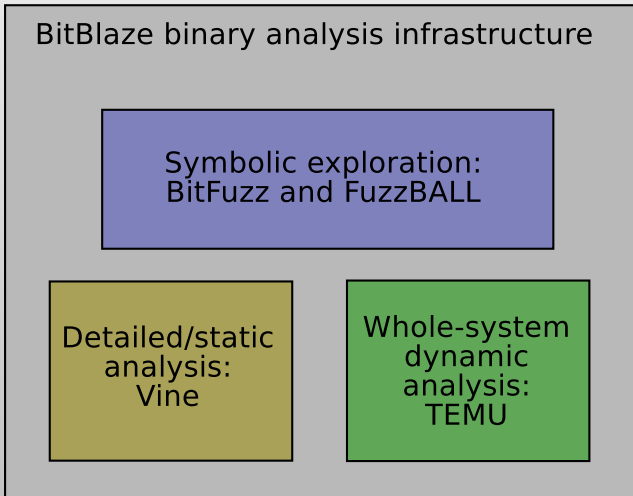
# Defense is challenging

- ❑ Software inevitably has bugs
- ❑ Attackers now have real incentives
  - Financial, national interest, ...
- ❑ Increasing sophistication and scale of attacks
- ❑ We need a new generation of defense techniques
  - Move beyond symptom-based and heuristic approaches

# The BitBlaze approach

- Use program semantics, focus on root causes
- Build a unified binary analysis platform for security
  - Leverage advances in program analysis, instrumentation, etc.
- Apply it to solve real-world security problems
  - I'll discuss just a few examples

# BitBlaze core components



# Outline

Core technique: symbolic reasoning

Binary-level bug-finding

Binary-level influence measurement

Real-time URL spam filtering

Strings and JavaScript vulnerabilities

## Basic idea

- ❑ Choose some of state (e.g., program or function input) to be symbolic:  
introduce variables for their values
- ❑ Computations on symbolic state  
produce formulas rather than concrete  
(e.g., integer) values
- ❑ Construct queries with these formulas,  
solve to answer questions about  
possible program behavior

## Why symbolic reasoning?

- + *Precise*: formulas can capture exact program behavior without approximation
- + *Complete solver*. (i.e. *decision procedure*) will always produce a correct solution without human help
- + *Flexibility*. Formulas independent of particular form of query



## Why not symbolic reasoning?

- Precise, but often not complete: don't prove that a given behavior can never happen
- Complete solver, but solution not guaranteed within reasonable space/time
- Flexibility, but may be be less efficient than more specialized approach



# Applications

Vulnerability signatures [Oakland'06,CSF'07] Protocol replay  
[CCS'06] Deviation discovery [USENIX'07] Patch-based exploit  
generation [Oakland'08] Modeling content sniffing [Oakland'09]  
Influence measurement [PLAS'09] Loop-extended SE  
[ISSTA'09] Protocol-level exploration [RAID'09] Kernel API  
exploration Decomposing crypto functions [CCS'10] Fixing  
under-tainting [NDSS'11] Protocol-model assisted SE [USENIX'11]  
JavaScript SE [Oakland'10] Static-guided test generation  
[ISSTA'11] Emulator verification [submitted]

# Applications

Vulnerability signatures [Oakland'06,CSF'07] Protocol replay

[CCS'06] Deviation discovery [USENIX'07] Patch-based exploit generation [Oakland'08] Modeling content sniffing [Oakland'09]

**Influence measurement [PLAS'09]** Loop-extended SE

[ISSTA'09] Protocol-level exploration [RAID'09] Kernel API exploration Decomposing crypto functions [CCS'10] Fixing

under-tainting [NDSS'11] Protocol-model assisted SE [USENIX'11]

**JavaScript SE [Oakland'10] Static-guided test generation**

**[ISSTA'11]** Emulator verification [submitted]

# Challenges of binary symbolic reasoning

## ❑ Instruction set complexity

- Rewrite to simpler intermediate language

## ❑ Variable-size memory accesses

- Lazy conversion with mixed-granularity storage

## ❑ No type distinction between integers and pointers

- Analyze symbolic expression structure

# Outline

Core technique: symbolic reasoning

**Binary-level bug-finding**

Binary-level influence measurement

Real-time URL spam filtering

Strings and JavaScript vulnerabilities

## Setting: vulnerability finding

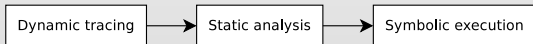
- Find exploitable bugs in software, before the bad guys do
- Many bugs found by independent researchers, without benefit of source code
- Example vulnerability type: buffer overflow
  - Incorrect or missing bounds check allows malicious input to overwrite other sensitive state
  - Despite extensive research, and some progress in practice, still a major bug category in C/C++ programs

# Static analysis

- ❑ Widely used at source-code level
- ❑ Can be *sound* (report all potential problems), at cost of false positives (*imprecision*)
- ❑ Challenge 1: more difficult at binary level
  - Soundness/precision tradeoff less favorable
- ❑ Challenge 2: developers have a low tolerance for false positives
  - Won't use a tool that wastes their time



# Combined static/dynamic approach



- Before static analysis, use dynamic traces to help where static binary analysis has trouble (e.g., indirect control flow)
- Design and optimize static analysis for binary-level challenges (e.g., variable identification, overlapping memory accesses)
- After static analysis, prioritize true positives by searching for test cases with symbolic execution

# Combined static/dynamic approach



- Before static analysis, use dynamic traces to help where static binary analysis has trouble (e.g., indirect control flow)
- Design and optimize static analysis for binary-level challenges (e.g., variable identification, overlapping memory accesses)
- After static analysis, prioritize true positives by searching for test cases with symbolic execution**

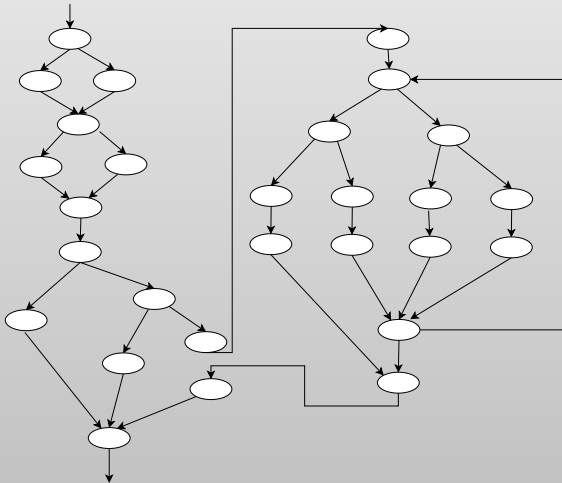
# Key challenge: guiding the search

- ❑ Increase the chances that the paths we explore will lead to a bug
  - Path must reach the code location of the bug
  - Program state at that location must trigger the bug
- ❑ Combination of two approaches:
  1. Data-flow slice and control-flow distance to direct paths toward a potential bug
  2. Explore patterns of loop body paths to cover cases likely to overflow

# Key challenge: guiding the search

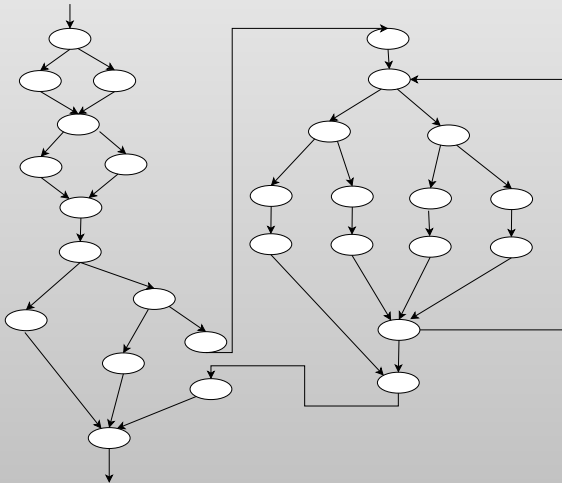
- ❑ Increase the chances that the paths we explore will lead to a bug
  - Path must reach the code location of the bug
  - Program state at that location must trigger the bug
- ❑ Combination of two approaches:
  1. Data-flow slice and control-flow distance to direct paths toward a potential bug
  2. Explore patterns of loop body paths to cover cases likely to overflow

# Guidance toward a bug

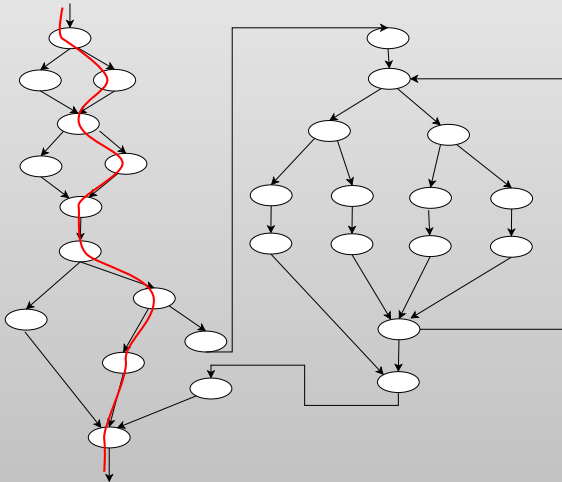




# Guidance toward a bug



# Guidance toward a bug

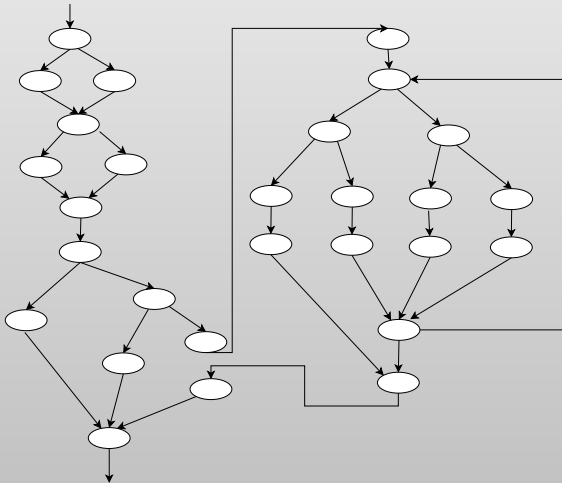




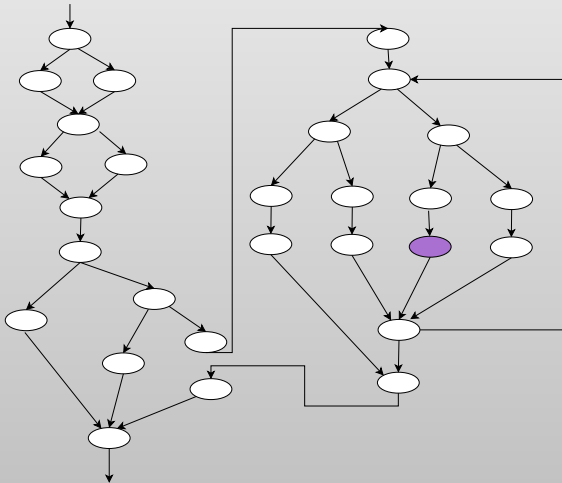




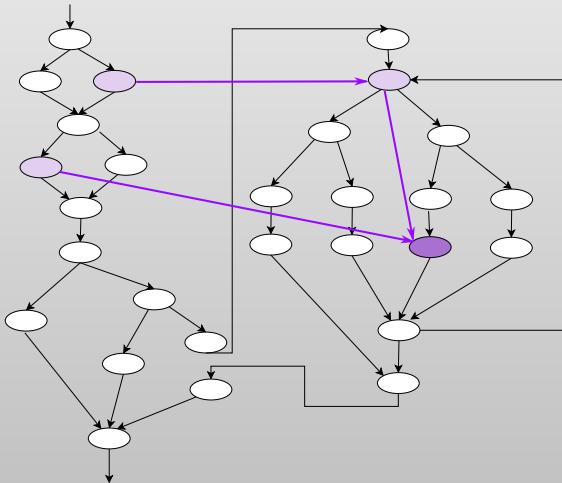
# Guidance toward a bug



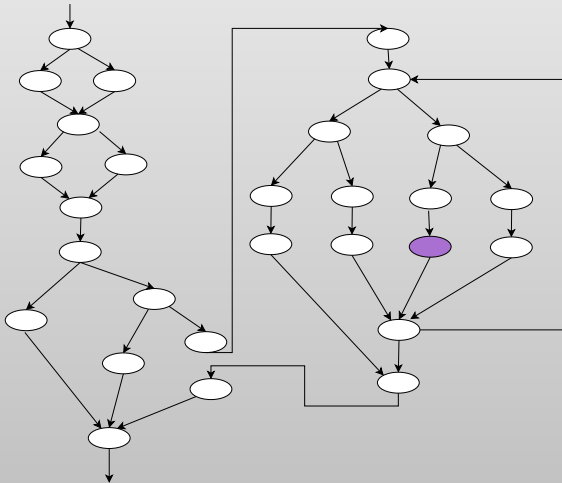
# Guidance toward a bug



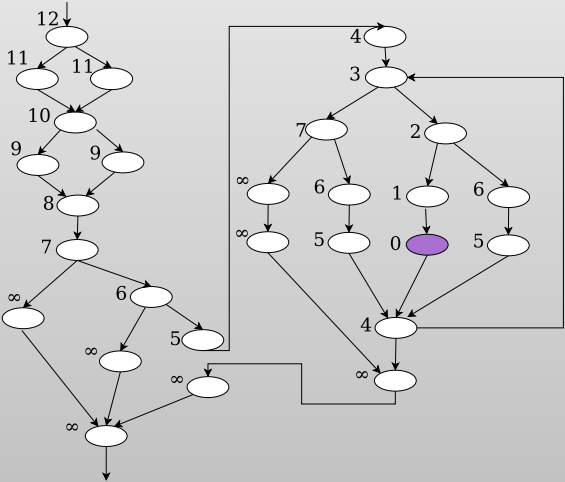
# Guidance toward a bug



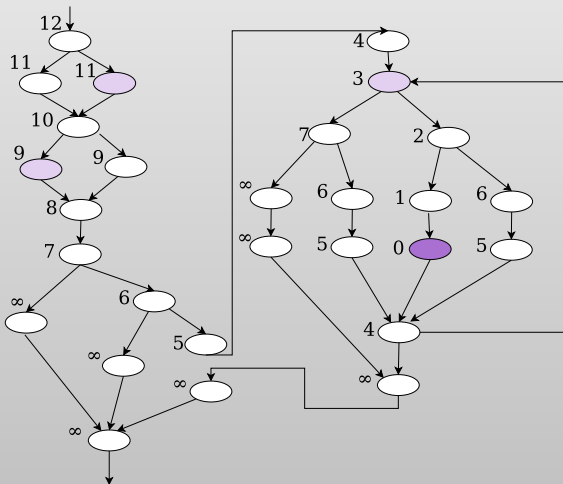
# Guidance toward a bug



# Guidance toward a bug



# Guidance toward a bug







## Sub-problem: control-flow distance

- An *interprocedural control-flow graph* has nodes for statements, and edges between statements and for calls and returns
- However, we can't use a regular graph distance measure (Dijkstra's algorithm), because of call and return matching
  - Exclude:  $f$  calls  $g$ ,  $g$  returns to  $h$
- Instead, new two-phase distance algorithm that first computes entry-to-exit distances bottom up, then adds unmatched returns and calls

# Guidance results

Benchmark	Unguided		Guided	
	Paths	Time (s)	Paths	Time (s)
BIND/b4	1	1.9	1	1.8
Sendmail/s5	3	19.0	3	22.9
BIND/b1	54	2.8	20	3.6
BIND/b2	137	13.3	72	25.1
BIND/b3	9	1.6	4	2.6
Sendmail/s2	16	2.9	9	97.0
Sendmail/s7	56	6.9	1	1.9
WU-FTPD/f1	309	8.1	11	1.1
WU-FTPD/f2	1455	65.8	11	1.4
WU-FTPD/f3	143	60.0	18	11.4
Sendmail/s5	T/O	> 21600.0	332	200.4
Sendmail/s6	T/O	> 21600.0	86	11.3
Sendmail/s1	T/O	> 21600.0	7297	7474.4
Sendmail/s3	T/O	> 21600.0	T/O	> 21600.0

# Guidance results

Benchmark	Unguided		Guided	
	Paths	Time (s)	Paths	Time (s)
→ BIND/b4	1	1.9	1	1.8
→ Sendmail/s5	3	19.0	3	22.9
BIND/b1	54	2.8	20	3.6
BIND/b2	137	13.3	72	25.1
BIND/b3	9	1.6	4	2.6
Sendmail/s2	16	2.9	9	97.0
Sendmail/s7	56	6.9	1	1.9
WU-FTPD/f1	309	8.1	11	1.1
WU-FTPD/f2	1455	65.8	11	1.4
WU-FTPD/f3	143	60.0	18	11.4
Sendmail/s5	T/O	> 21600.0	332	200.4
Sendmail/s6	T/O	> 21600.0	86	11.3
Sendmail/s1	T/O	> 21600.0	7297	7474.4
Sendmail/s3	T/O	> 21600.0	T/O	> 21600.0

# Guidance results

Benchmark	Unguided		Guided	
	Paths	Time (s)	Paths	Time (s)
BIND/b4	1	1.9	1	1.8
Sendmail/s5	3	19.0	3	22.9
→ BIND/b1	54	2.8	20	3.6
→ BIND/b2	137	13.3	72	25.1
→ BIND/b3	9	1.6	4	2.6
→ Sendmail/s2	16	2.9	9	97.0
Sendmail/s7	56	6.9	1	1.9
WU-FTPD/f1	309	8.1	11	1.1
WU-FTPD/f2	1455	65.8	11	1.4
WU-FTPD/f3	143	60.0	18	11.4
Sendmail/s5	T/O	> 21600.0	332	200.4
Sendmail/s6	T/O	> 21600.0	86	11.3
Sendmail/s1	T/O	> 21600.0	7297	7474.4
Sendmail/s3	T/O	> 21600.0	T/O	> 21600.0

# Guidance results

Benchmark	Unguided		Guided	
	Paths	Time (s)	Paths	Time (s)
BIND/b4	1	1.9	1	1.8
Sendmail/s5	3	19.0	3	22.9
BIND/b1	54	2.8	20	3.6
BIND/b2	137	13.3	72	25.1
BIND/b3	9	1.6	4	2.6
Sendmail/s2	16	2.9	9	97.0
→ Sendmail/s7	56	6.9	1	1.9
→ WU-FTP/f1	309	8.1	11	1.1
→ WU-FTP/f2	1455	65.8	11	1.4
→ WU-FTP/f3	143	60.0	18	11.4
Sendmail/s5	T/O	> 21600.0	332	200.4
Sendmail/s6	T/O	> 21600.0	86	11.3
Sendmail/s1	T/O	> 21600.0	7297	7474.4
Sendmail/s3	T/O	> 21600.0	T/O	> 21600.0

# Guidance results

Benchmark	Unguided		Guided	
	Paths	Time (s)	Paths	Time (s)
BIND/b4	1	1.9	1	1.8
Sendmail/s5	3	19.0	3	22.9
BIND/b1	54	2.8	20	3.6
BIND/b2	137	13.3	72	25.1
BIND/b3	9	1.6	4	2.6
Sendmail/s2	16	2.9	9	97.0
Sendmail/s7	56	6.9	1	1.9
WU-FTPD/f1	309	8.1	11	1.1
WU-FTPD/f2	1455	65.8	11	1.4
WU-FTPD/f3	143	60.0	18	11.4
→ Sendmail/s5	T/O	> 21600.0	332	200.4
→ Sendmail/s6	T/O	> 21600.0	86	11.3
Sendmail/s1	T/O	> 21600.0	7297	7474.4
Sendmail/s3	T/O	> 21600.0	T/O	> 21600.0

# Guidance results

Benchmark	Unguided		Guided	
	Paths	Time (s)	Paths	Time (s)
BIND/b4	1	1.9	1	1.8
Sendmail/s5	3	19.0	3	22.9
BIND/b1	54	2.8	20	3.6
BIND/b2	137	13.3	72	25.1
BIND/b3	9	1.6	4	2.6
Sendmail/s2	16	2.9	9	97.0
Sendmail/s7	56	6.9	1	1.9
WU-FTPD/f1	309	8.1	11	1.1
WU-FTPD/f2	1455	65.8	11	1.4
WU-FTPD/f3	143	60.0	18	11.4
Sendmail/s5	T/O	> 21600.0	332	200.4
Sendmail/s6	T/O	> 21600.0	86	11.3
→ Sendmail/s1	T/O	> 21600.0	7297	7474.4
Sendmail/s3	T/O	> 21600.0	T/O	> 21600.0



# Guidance results

Benchmark	Unguided		Guided	
	Paths	Time (s)	Paths	Time (s)
BIND/b4	1	1.9	1	1.8
Sendmail/s5	3	19.0	3	22.9
BIND/b1	54	2.8	20	3.6
BIND/b2	137	13.3	72	25.1
BIND/b3	9	1.6	4	2.6
Sendmail/s2	16	2.9	9	97.0
Sendmail/s7	56	6.9	1	1.9
WU-FTPD/f1	309	8.1	11	1.1
WU-FTPD/f2	1455	65.8	11	1.4
WU-FTPD/f3	143	60.0	18	11.4
Sendmail/s5	T/O	> 21600.0	332	200.4
Sendmail/s6	T/O	> 21600.0	86	11.3
Sendmail/s1	T/O	> 21600.0	7297	7474.4
→ Sendmail/s3	T/O	> 21600.0	T/O	> 21600.0

# Outline

Core technique: symbolic reasoning

Binary-level bug-finding

Binary-level influence measurement

Real-time URL spam filtering

Strings and JavaScript vulnerabilities

## Due and undue influence

- ❑ How much influence should network inputs have on a program?
- ❑ For instance, on an indirect jump target
  - Some influence → select a legal behavior
  - Too much influence → control flow hijacking attack

# High and low influence examples

```
void (*func_ptr)(void);  
func_ptr = untrusted_input();  
(*func_ptr)();
```

---

```
void (*func_ptr)(void);  
switch (untrusted_input()) {  
    case CMD_OPEN: func_ptr = &open_file;  
    case CMD_READ: func_ptr = &read_file;  
    default:        func_ptr = &error;  
}  
(*func_ptr)();
```

## Channel capacity as influence



- For a given variable, how many values can an attacker produce?
- Influence =  $\log_2(\# \text{ values})$
- Special case of channel capacity from information theory

## Scalability and precision

- ❑ Want to analyze large (e.g., commercial) software
- ❑ Want results with no error
- ❑ Our goal: improved trade-off points between these ideals

# Problem statement

## Given:

- A deterministic program with designated inputs
- An output variable

Question: how many values of the output are possible, given different inputs?

## Program to formula example

```
/* Convert low 4 bits of integer to hex */
char tohex(int i) {
    int low = i & 0xf;
    char v;
    if (low < 10)
        v = '0' + low;
    else
        v = 'a' + (low - 10);
    return v;
}
```

**Dynamic:**  $(i \& 15) < 10 \wedge (v = 48 + (i \& 15))$



## Program to formula example

```
/* Convert low 4 bits of integer to hex */
char tohex(int i) {
    int low = i & 0xf;
    char v;
    if (low < 10)
        v = '0' + low;
    else
        v = 'a' + (low - 10);
    return v;
}
```

**Static:**  $((i \& 15) < 10 \wedge (v = 48 + (i \& 15))) \vee$   
 $((i \& 15) \geq 10 \wedge (v = 97 + (i \& 15) - 10))$

# Query techniques



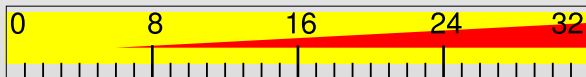
- ❑ Point-by-point exhaustion
- ❑ Range exclusion
- ❑ Random output sampling
- ❑ Probabilistic model counting

# Query techniques



- ❑ Point-by-point exhaustion
- ❑ Range exclusion
- ❑ **Random output sampling**
- ❑ Probabilistic model counting

# Query techniques



- Point-by-point exhaustion
- Range exclusion
- Random output sampling
- Probabilistic model counting

## Point-by-point exhaustion

- Is  $v = f(i)$  satisfiable?
- Suppose it is, by  $v_1 = f(i_1)$
- Is  $v = f(i) \wedge v \neq v_1$  satisfiable?
- ...
- We repeat up to at most  $2^6 = 64$  distinct outputs, so every bound up to 6 bits is exact

## Range exclusion

- ❑ Is  $v = f(i) \wedge (a \leq v \leq b)$  satisfiable?
- ❑ If not, a whole range is excluded
- ❑ If so, can subdivide
- ❑ We also use this with binary search to find the minimum and maximum outputs

## Random output sampling

- ❑ Pick  $v_r$  at random, and check if  $v_r = f(i)$  is satisfiable
- ❑ By default, our tool uses 20 samples, and computes a 95% confidence interval

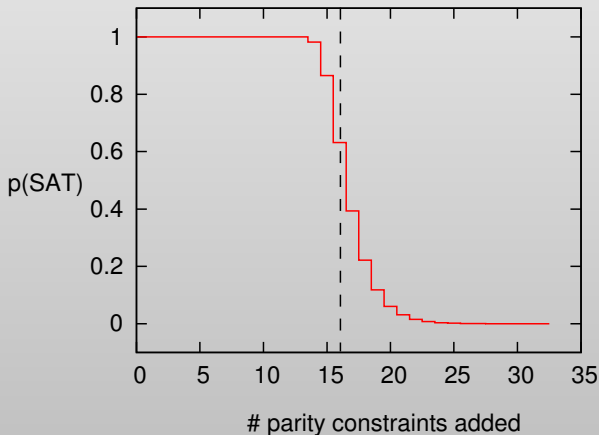
## Probabilistic model counting

- Use XOR streamlining [GSS06] to probabilistically reduce #SAT to SAT
- Analogy: counting audience members
- Random parity constraints over enough bits are effectively independent
- Perform repeated experiments with different numbers of constraints



# Probabilistic model counting

Choose # of constraints so that  $p(\text{SAT}) \approx 0.5$








# Identity function

$$v = i$$

Low	High	Sample	#SAT	Actual
6.04	32.0	[31.8, 32.0]	32.0	32

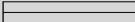


-  Feasible Point
-  Infeasible Range
-  ~100% (Probabilistic)
-  ~50% (Probabilistic)
-  < 5% (Probabilistic)


# tohex

```
sprintf(&v, "%x", i & 0xf)
```

Static:

Low	High	Sample	#SAT	Actual	
4.00	4.00	N/A	N/A	4	

Dynamic:

Low	High	Sample	#SAT	Actual	
3.32	3.32	N/A	N/A	$\log_2 10$	
2.58	2.58	N/A	N/A	$\log_2 6$	

## Mix and duplicate

$$f(x \circ y) = (x \oplus y) \circ (x \oplus y)$$

$$f(0x00000042) = 0x00420042$$

$$f(0x02461111) = 0x13571357$$

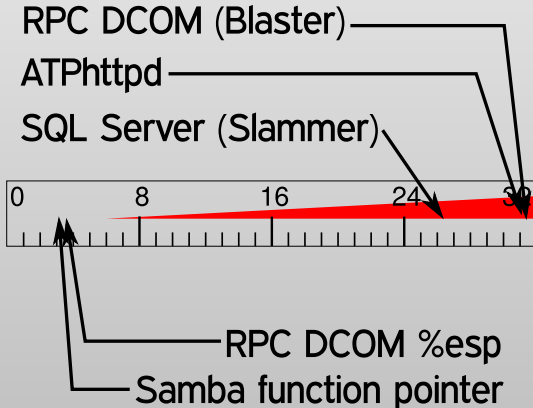
$$f(0xcafebebe) = 0x74407440$$

Low	High	Sample	#SAT	Actual
6.04	32.0	[0.0, 28.6]	15.8	16






# Results summary

Goal: distinguish attacks from false positives



## Confirming attacks

- ❑ Vulnerable Windows and Linux binaries
- ❑ Real attacks all have high influence, at least 26 bits

Program	High	Sample	#SAT	Value Set
RPC DCOM	32.0	[31.8, 32.0]	30.4	
SQL Server	30.9	[26.7, 28.3]	26.6	
ATPhttpd	32.0	[31.8, 32.0]	31.0	

## Reveal false positives

- Examples cause taint analysis warnings
- Measured influence exactly, less than 5 bits

Program	Low	High	Value Set
RPC %esp	3.81	3.81	<input type="text"/>
Samba func. ptr	3.32	3.32	<input type="text"/>

# Directions for improving solving

- Further targeted query strategies
  - E.g., two-bit patterns [Meng & Smith, PLAS'11]
- Refined strategy for choosing number of parity constraints
- Interface with off-the-shelf #SAT solvers
  - Question: how to restrict counting to output bits?



# Outline

Core technique: symbolic reasoning

Binary-level bug-finding

Binary-level influence measurement

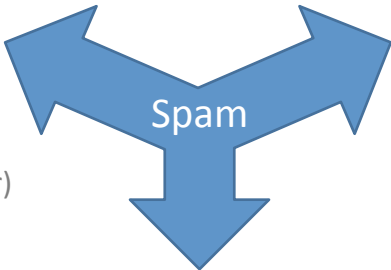
Real-time URL spam filtering

Strings and JavaScript vulnerabilities

# Motivation



**Social Networks**  
(Facebook, Twitter)



**Blogs, Services**  
(Blogger, Yelp)



**Web Mail**  
(Gmail, Live Mail)

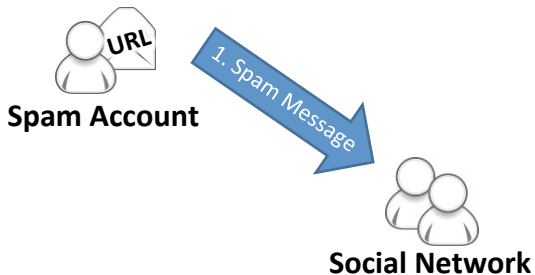
# Motivation

- Existing solutions:
  - Blacklists
  - Service-specific, account heuristics
- Develop new spam filter service:
  - Filter spam: scams, phishing, malware
  - Real-time, fine-grained, generalizable

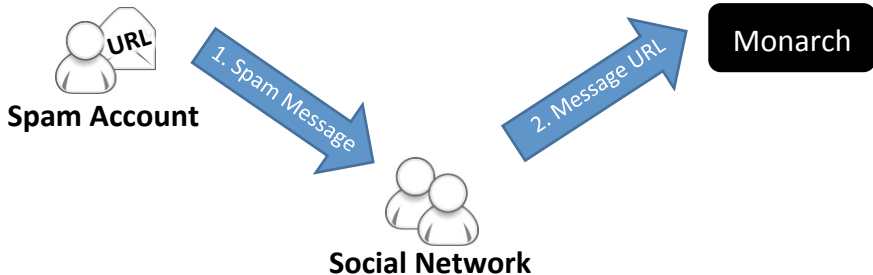
# Overview

- Our system – **Monarch**:
  - Accepts millions of URLs from web service
  - Crawls, labels each URL in real-time
- Spam Classification
  - Decision based on URL content, page behavior, hosting
  - Large-scale; distributed collection, classification
- Implemented as a cloud service

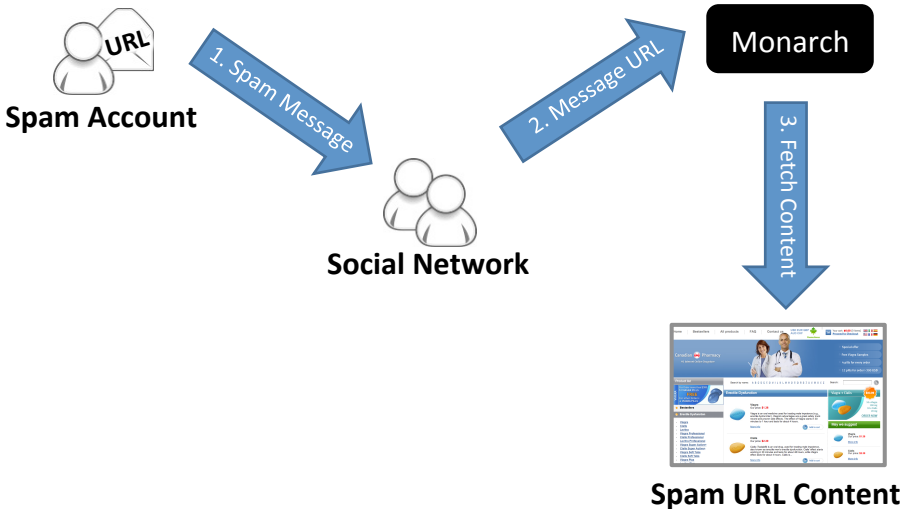
# Monarch in Action



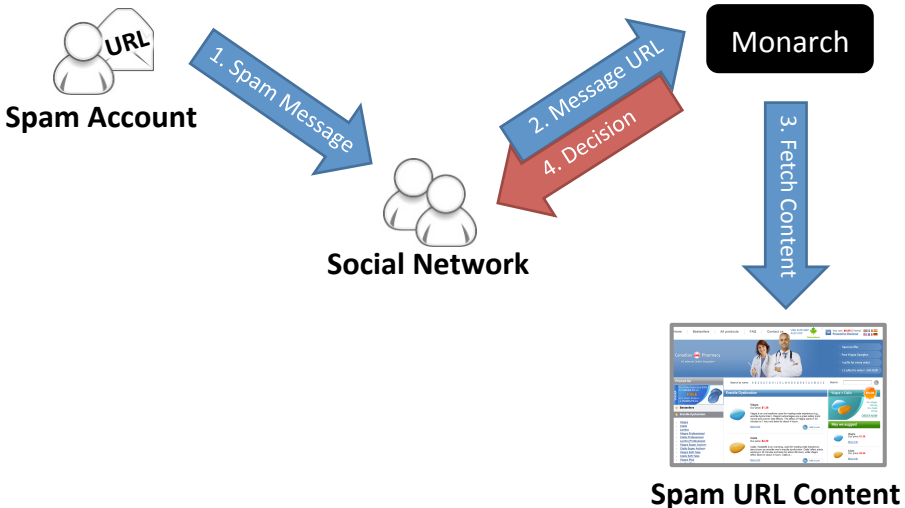
# Monarch in Action



# Monarch in Action

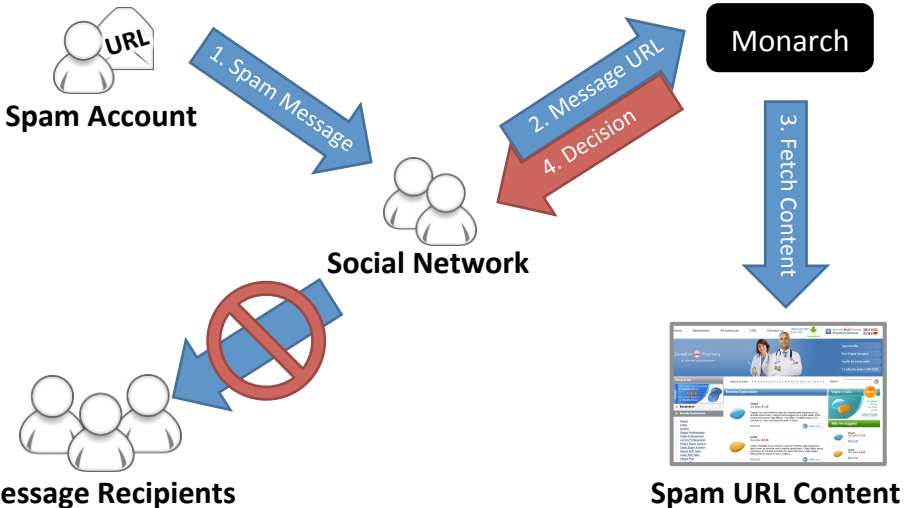


# Monarch in Action





# Monarch in Action



# Challenges

Accuracy

Real-Time

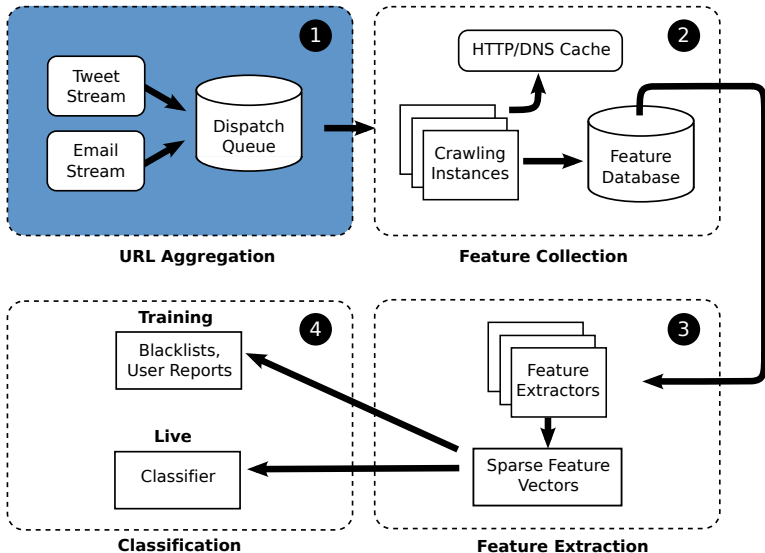
Scalability

Tolerant to Feature Evolution

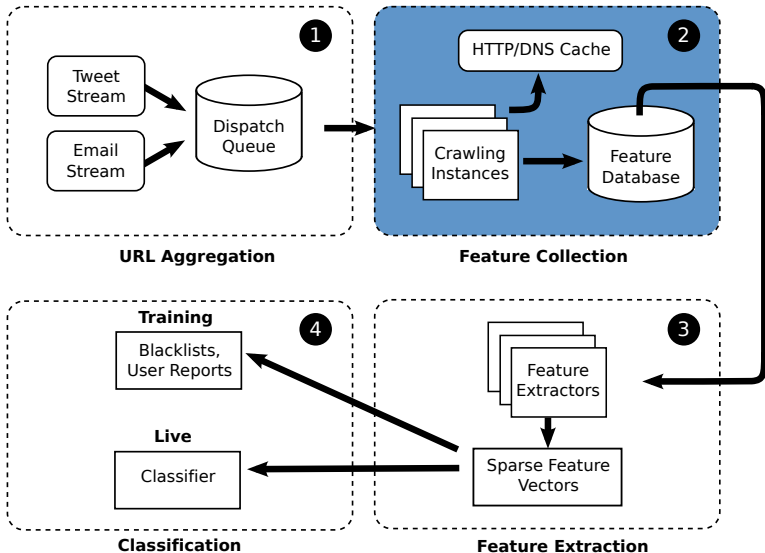
# Outline

- Architecture
- Results & Performance
- Limitations
- Conclusion

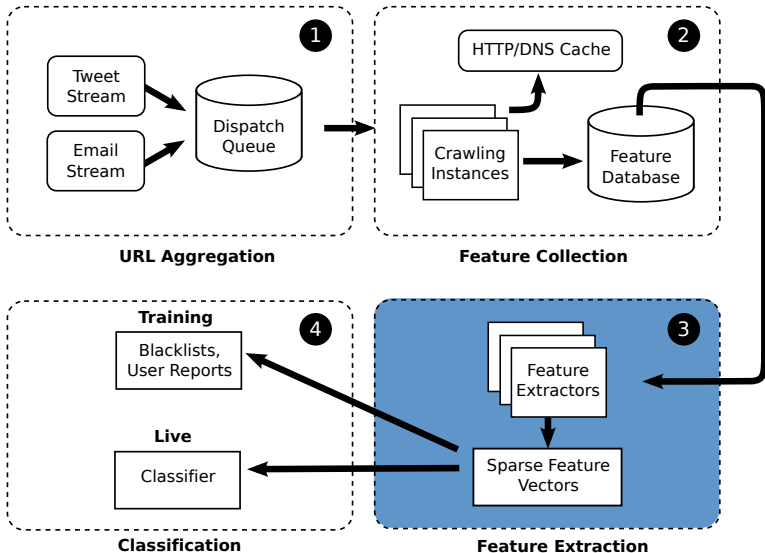
# System Architecture



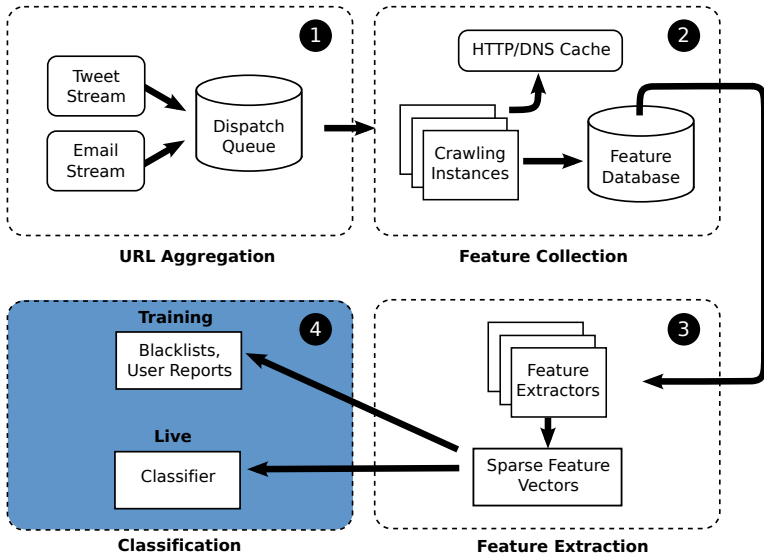
# System Architecture



# System Architecture



# System Architecture



# URL Aggregation

Source	Sample Size
Spam email URLs	1.25 million
Blacklisted Twitter URLs	567,000
Non-spam Twitter URLs	9 million

Collection period: 9/8/2010 – 10/29/2010



# Feature Collection

- **High Fidelity Browser**
- **Navigation**
  - Lexical features of URLs (length, subdomains)
  - Obfuscation (directory operations, nested encoding)
- **Hosting**
  - IP/ASN
  - A, NS, MX records
  - Country, city if available

# Feature Collection

- **Content**

- Common HTML templates, keywords
- Search engine optimization
- Content of request, response headers

- **Behavior**

- Prevent navigating away
- Pop-up windows
- Plugin, JavaScript redirects

# Classification

$$f(\vec{w}) = \sum_{i=1}^n \log(1 + \exp[-y_i(\vec{x}_i \cdot \vec{w}_i)])$$

---

- **Distributed Logistic Regression**
  - Data overload for single machine

# Classification

$$f(\vec{w}) = \sum_{i=1}^n \log(1 + \exp[-y_i(\vec{x}_i \cdot \vec{w}_i)]) + \lambda \|\vec{w}\|_1.$$

---

- **Distributed Logistic Regression**
  - Data overload for single machine
- **L1-regularization**
  - Reduces feature space, over-fitting
  - 50 million features -> 100,000 features

# Implementation

- System implemented as a cloud service on Amazon EC2
  - **Aggregation:** 1 machine
  - **Feature Collection:** 20 machines
    - Firefox, extension + modified source
  - **Classification & Feature Extraction:** 50 machines
    - Hadoop - Spark, Mesos
- Straightforward to scale the architecture

# Result Overview

- High-level summary:
  - Performance
  - Overall accuracy
  - Highlight important features
  - Feature evolution
  - Spam independence between services

# Performance

- Rate: 638,000 URLs/day
  - Cost: \$1,600/mo
- Process time: 5.54 sec
  - Network delay: 5.46 sec
- Can scale to 15 million URLs/day
  - Estimated \$22,000/mo

# Measuring Accuracy

- Dataset: 12 million URLs (<2 million spam)
  - Sample 500K spam (half tweets, half email)
  - Sample 500K non-spam
- Training, Testing
  - 5-fold validation
  - Vary training folds non-spam:spam ratio
  - Test fold equal parts spam, non-spam



# Overall Accuracy

Training Ratio	Accuracy	False Positive Rate	False Negative Rate
1:1	94%	4.23%	7.5%
4:1	91%	0.87%	17.6%
10:1	87%	0.29%	26.5%



**Correctly labeled samples**



**Non-spam labeled as spam**



**Spam labeled as non-spam**

# Overall Accuracy

Training Ratio	Accuracy	False Positive Rate	False Negative Rate
1:1	94%	4.23%	7.5%
4:1	91%	0.87%	17.6%
10:1	87%	0.29%	26.5%



Correctly labeled samples

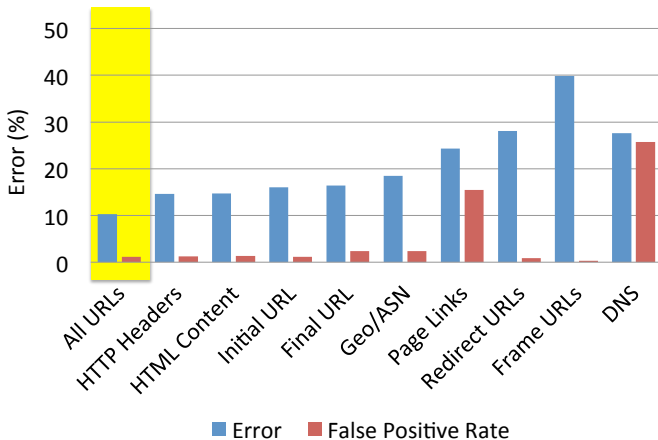


Non-spam labeled as spam



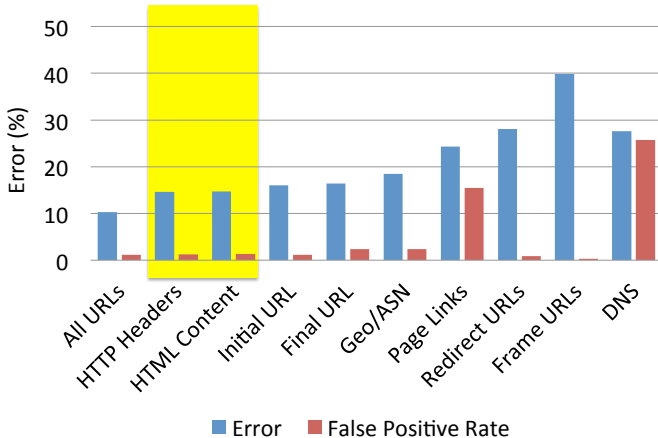
Spam labeled as non-spam

# Error by Feature



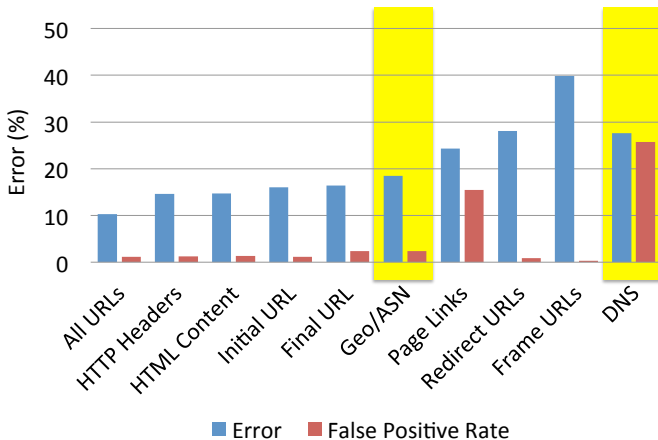
**Error = 1 - Accuracy**

# Error by Feature



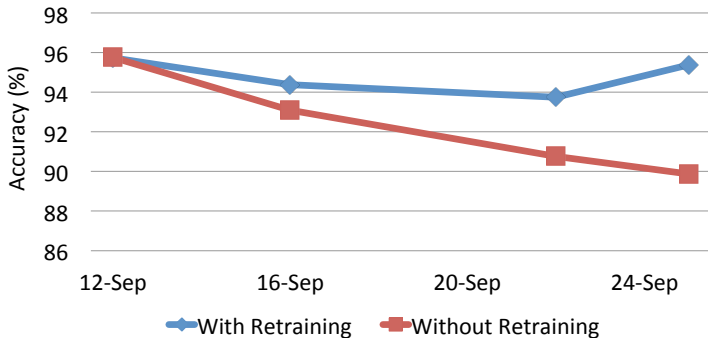
**Error = 1 - Accuracy**

# Error by Feature



**Error = 1 - Accuracy**

# Feature Evolution – Retraining Required



# Spam Independence

- Unexpected result: Twitter, email spam qualitatively different

Training Set	Testing Set	Accuracy	False Negatives
<b>Twitter</b>	<b>Twitter</b>	<b>94%</b>	<b>22%</b>
Twitter	Email	81%	88%
Email	Twitter	80%	99%
<b>Email</b>	<b>Email</b>	<b>99%</b>	<b>4%</b>

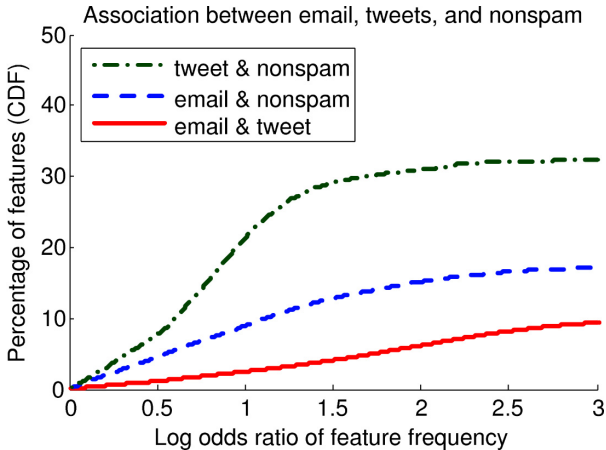
# Spam Independence

- Unexpected result: Twitter, email spam qualitatively different

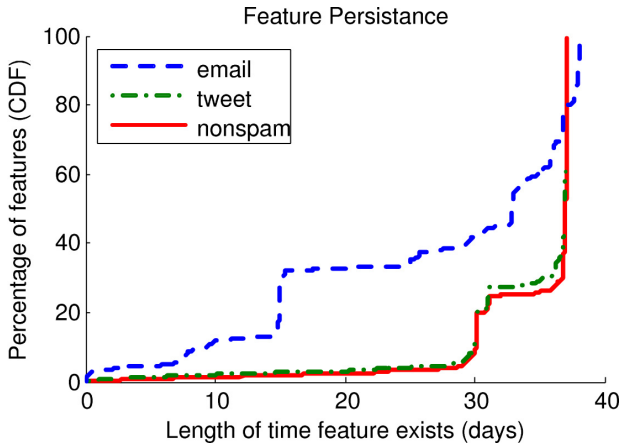
Training Set	Testing Set	Accuracy	False Negatives
Twitter	Twitter	94%	22%
<b>Twitter</b>	<b>Email</b>	<b>81%</b>	<b>88%</b>
<b>Email</b>	<b>Twitter</b>	<b>80%</b>	<b>99%</b>
Email	Email	99%	4%



# Distinct Email, Twitter Features



# Email Features Shorter Lived



# Limitations

- Adversarial Machine Learning
  - We provide oracle to spammers
  - Can adversaries tweak content until passing?
- Time-based Evasion
  - Change content after URL submitted for verification
- Crawler Fingerprinting
  - Identify IP space of Monarch, fingerprint Monarch browser client
  - Dual-personality DNS, page behavior

# Outline

Core technique: symbolic reasoning

Binary-level bug-finding

Binary-level influence measurement

Real-time URL spam filtering

Strings and JavaScript vulnerabilities

# Example attack: gadget overwrite

The screenshot shows a Firefox browser window with the address bar displaying `http://www.google.com/ig`. The page content is the iGoogle homepage, but the TVGuide.com gadget has been overwritten. The gadget's title bar reads "TVGuide.com". The main content of the gadget is:

**TVGuide iGoogle gadget now *integrated* with iGoogle!**  
Please login to view the TVGuide listings.

Sign in to iGoogle with your  
**Google Account**

Email:

Password:

Stay signed in

[Can't access your account?](#)

Don't have a Google Account?  
[Create an account now](#)

## Example attack: explanation

- Cross-site scripting can exist entirely in client-side JavaScript
- Unsanitized data passed to HTML creation (`document.write`) or `eval`
- In the example, a malicious link injects code into the TVGuide gadget, turning it into a phishing vector

## What's new here?

- Source/sink problem, somewhat like SQL injection or server-side XSS, but:
  - JS code takes many kinds of inputs as unstructured strings, requiring custom parsing
  - Sanitization is not standardized, and often application-specific

→ More difficult challenges for string reasoning

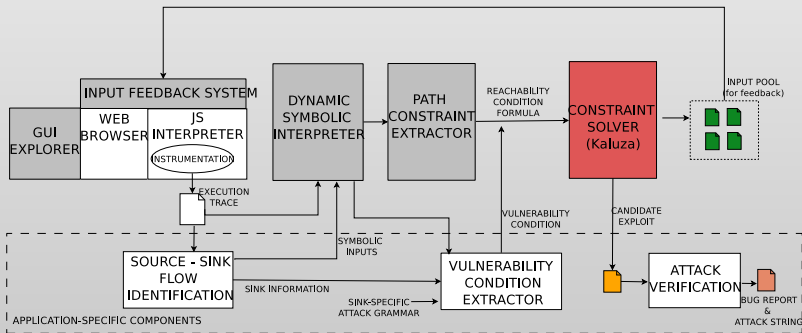
# Exploration overview

Two kinds of exploration:

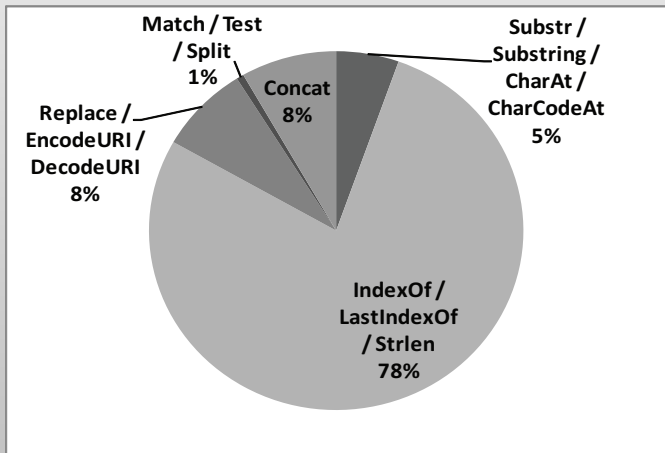
- ❑ Event space: GUI actions such as clicking check-boxes or links
- ❑ Value space: contents of form, message, and URL inputs
  - Explore new program paths
  - Check whether sanitization is sufficient (compare to attack grammar)



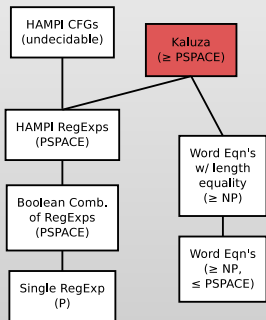
# Kudzu system overview



# Usage of string operations



# Expressiveness

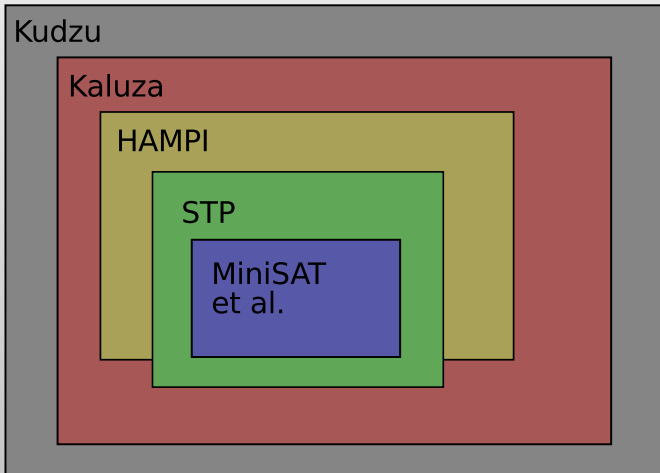


(Complexities are for unbounded variants)

- Regular expression membership
- Arbitrary concatenation (word equations)
- String length function

Can also mix in boolean and (31-bit) integer constraints

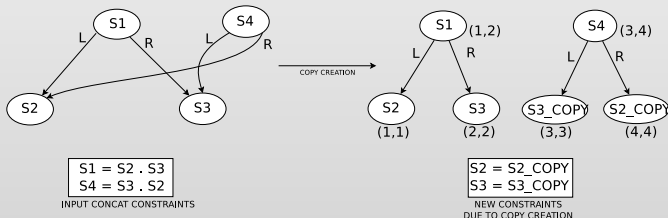
# Nested architecture



## Approach overview

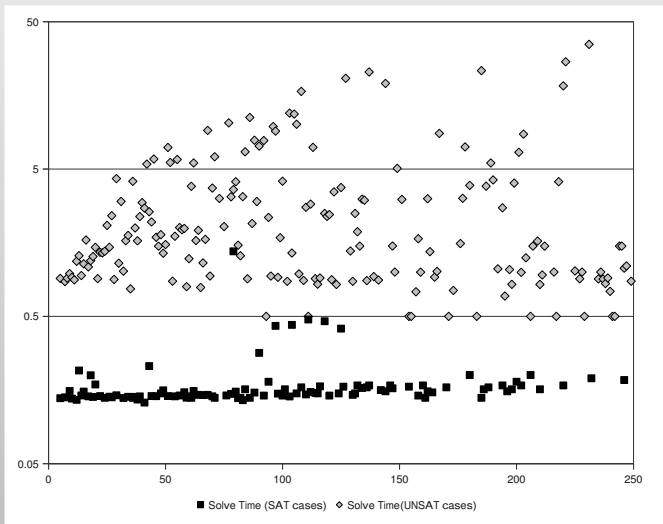
- ❑ Flatten concatenations to a linear sequence
- ❑ Abstract to length constraints
- ❑ For each length assignment:
  - Expand regexps (HAMPI code)
  - Combine in single bitvector query
  - bitvector SAT  $\rightarrow$  string SAT
- ❑ Exhausted lengths  $\rightarrow$  string UNSAT

# Approach details



- Real JavaScript “regexes” are more complex than textbook ones
- Regex lengths  $\rightarrow$  ultimately periodic set
- Translate `replace` with fixed number of occurrences

# Kaluza performance results



## Overall results

- Tested 5 AJAX applications and 13 iGoogle gadgets (all live)
- Event and value space exploration both contribute to coverage
  - But some code and events not yet covered
- Found vulnerabilities in 11 apps, including 2 missed by our previous taint-directed fuzzer



## Summary, and for more info

- ❑ Symbolic exploration and reasoning enable a wide variety of security tools
- ❑ Many security problems need lots of computing, but are naturally parallelizable
- ❑ <http://bitblaze.cs.berkeley.edu/>
- ❑ <http://bit.ly/muahjS>

Thank you

# Backup slides

## Web browser content sniffing

- An HTTP response contains a content type header
  - E.g., `text/html` or `image/png`
- But sometimes (~1%) the content type is missing or invalid
- Thus browsers sometimes attempt to sniff (guess) the type from the content or URL

## When content sniffing goes bad

- ❑ Content type matters because it affects privilege
  - Some types of content (HTML, Flash) can contain code
- ❑ An unexpected upgrade can allow an untrusted user to inject JavaScript
  - I.e., a kind of cross-site scripting (XSS)
- ❑ Usually a mismatch between the browser and another filter

## HotCRP attack example

- ❑ Conference site allows authors to upload PostScript papers
- ❑ What if the site accepts this file as PS, but the reviewer's browser considers it HTML?

```
%!PS-Adobe
%%Creator: <script>submitReview("A+");
...
```

- ❑ Your paper gets accepted :-)

## Modeling content sniffing

- To understand such attacks, we want a formal model of the sniffer's behavior
- E.g.,  $M^H(c) = \text{true}$  if the file contents  $c$  are sniffed as HTML
- Boolean combinations correspond to possible mismatch attacks
  - $M_1^P(c) \wedge M_2^H(c)$

## Model extraction

- The content-sniffing strategies of closed-source browsers are often un- or under-documented
  - We look at IE 7, Safari 3.1
- Extract from the binary using white-box exploration (symbolic execution)
- Model is a disjunction of path conditions from accepting paths



## Abstracting string functions

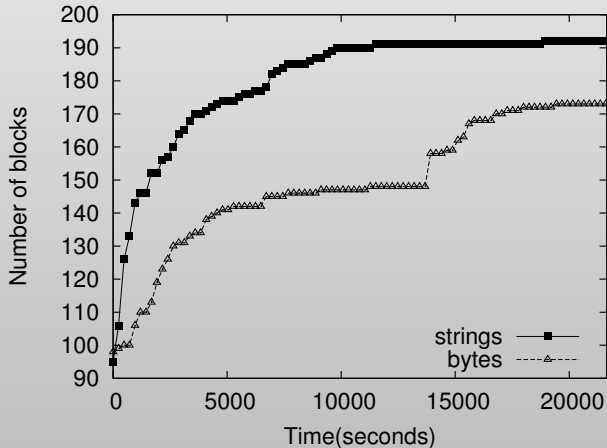
- ❑ Sniffing code makes heavy use of string routines
- ❑ Reason about their semantics, not their implementation
  - + Summarize multiple paths
  - + Skip implementation details
  - + Take advantage of specialized solvers (future)

## Translating string functions

1. Recognize over 100 binary-level functions (mostly documented)
2. Canonicalize to 14 semantic classes
3. Express in terms of a core constraint language
4. Reduce core constraints to STP bit vectors

# Exploration advantage of strings

Block coverage for Safari:



## Summary of attacks found

- ❏ Tool finds attacks to upgrade 6 content types each in IE and Safari to HTML.
  - But which pass a common server-side filter
  - Wikipedia has a more complex filter, but it can also be bypassed
- ❏ Automatically generated PS → HTML example:

```
%!t?HPTw\n0tKoCg1D<HeadswssssRsD
```

## Happy ending: safe sniffing

- Our models can be used to create matching server-side filters
- We propose client-side design principles for safe sniffing
  - Avoid privilege escalation
  - Prefix-disjoint signatures
- Adopted by IE 8 (partial), Chrome, and HTML 5

## Guidance (2): loop patterns

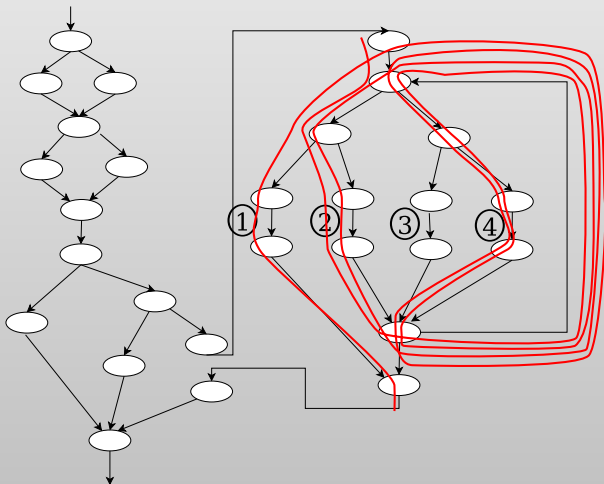
- Observation: triggering a loop-related bug often requires a specific execution pattern within a loop

```
char buf[20], *p = buf;
int mode = 0;
while (p >= buf) {
    switch (read_char()) {
        case 'a':
            if (mode == 1) {
                *p++ = 'x'; mode = 0; // path 1
            } else { // path 2
                } break;
        case 'b':
            mode = !mode; break; // path 3
        default: p = max(buf, p--); // path 4
    }
}
```

## Guidance (2): loop patterns

- ❑ Approach: cover a variety of patterns during symbolic execution
  - Don't try to find the right pattern statically
- ❑ Statically number paths through a loop body
- ❑ Try patterns in inverse relation to their length
- ❑ Interleave use of patterns with discovering which paths are feasible

## Guidance (2): loop patterns



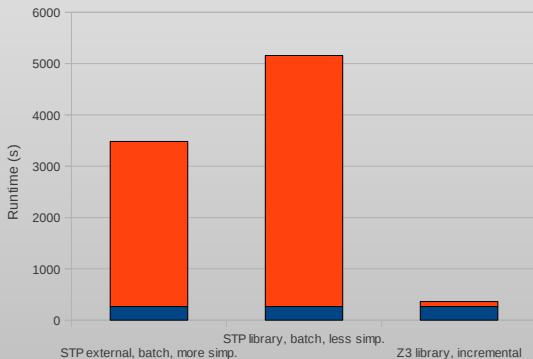


## What do our formulas look like?

- The key theory is fixed-size bit-vectors, representing machine integers
  - Exact treatment of overflow, signs, etc. important for binaries
- Could use arrays for general memory, lookup tables, but usually don't
  - Instead, fix memory layout to be concrete (or unconstrained symbolic)
- Usually easy to solve, whether SAT or UNSAT

# Solver performance


For easy formulas, mundane changes matter (sample of 84355 formulas, not a general tool comparison)



# Checked/bounded copy

```
v = i & 0x0f
```

```
v = 0; if (i < 16) v = i
```

Low	High	Sample	#SAT	Actual	
4.00	4.00	N/A	N/A	4	

# Multiplication and division

$$v = i * 2$$

Low	High	Sample	#SAT	Actual
6.58	32.0	[30.4, 31.6]	31.5	31



$$v = i / 2$$

Low	High	Sample	#SAT	Actual
6.58	31.0	[30.8, 31.0]	31.7	31

