



blackhat
USA 2012

libinjection

a C library for SQLi detection and generation
through lexical analysis of real world attacks

Nick Galbreath @ngalbreath nickg@client9.com

Wednesday July 25, 2012 2:45PM
Augustus I+II, Caesar's Palace, Las Vegas



blackhat
USA 2012

The latest version of this presentation:

<http://slidesha.re/>

OBch5k

whoami

- ▶ Director of Engineering @ Etsy
- ▶ Enterprise, Fraud, Security, Email, Fun
- ▶ stringencoders:
 - ▶ C library for string processing
 - ▶ used by every ad server in the world
 - ▶ used in Chrome browser
 - ▶ <http://code.google.com/p/stringencoders>

The Next 14 Minutes

- ▶ Why is detecting SQLi hard
- ▶ The algorithm behind libinjection
- ▶ The results
- ▶ Next Steps

Detecting SQLi
from User Input
is a Hard Problem

It's Easy to Get Started with Regular Expressions!

s/UNION\s+(ALL)?/i

- ▶ At least two open source WAF use regular expressions.
- ▶ Failure cases in closed-source WAFs also indicate regexp.

SQL IS HUGE

- ▶ Turing Complete! (sorta)
- ▶ 1992 SQL Spec: bit.ly/10fmhZ
 - ▶ 625 pages of plain text
- ▶ 2003 SQL Spec: bit.ly/0B5vfW
 - ▶ 128 pages of pure BNF
- ▶ No one implements exactly
- ▶ Everyone has extensions, exceptions, bugs

It's more complicated than you think.

- ▶ Recursive commenting rules
- ▶ A single *number* can't be done in a single regexp.
- ▶ Really Loosely Typed
- ▶ String rules - OMFG. You think you know but you have no idea.
- ▶ Come see my talk at DEFCON this Friday at...
4:20 pm

Guns and Butter

- ▶ In 2005, right here at Black Hat, Hanson and Patterson presented:
Guns and Butter: Towards Formal Axioms of Validation (<http://bit.ly/OBe7mJ>)
- ▶ *...formally proved that for any regex validator, we could construct either a safe query which would be flagged as dangerous, or a dangerous query which would be flagged as correct.*
- ▶ *(summary from libdejector documentation)*

Existing WAFs

- ▶ Visual inspection shows bugs
- ▶ Don't see very much in testing
- ▶ Don't see much or any false positive testing
- ▶ Closed source WAF have zero accountability (e.g. there is no formal disclosure of what they detect or not, and how they do it)

CAN WE DO BETTER?

Nick Galbreath



@ngalbreath

libinjection

Key Insight

- ▶ A SQLi attack must be parsed as SQL with the original query.
- ▶ "Is it a SQLi attack?" becomes "Could it be a SQL snippet?"

Only 3 Contexts

User input is only "injected" into SQL in three ways:

- ▶ "As-Is"
- ▶ Inside a single quoted string
- ▶ Inside a double quoted string

(I suppose another would be inside a comment, but we can't do everything)

Identification of SQL snippets without context is hard

- ▶ 1-917-660-3400 my phone number or an arithmetic expression?
- ▶ @ngalbreath my twitter account or a SQL variable?

Existing SQL Parsers

- ▶ Only parse their flavor of SQL
- ▶ Not well designed to handle snippets
- ▶ Hard to extend
- ▶ Worried about correctness

... so I wrote my own!

Tokenization

- ▶ Converts input into a stream of tokens
- ▶ Uses "master list" of keywords and functions across all databases.
- ▶ Handles comments, string, literals, weirdos.

```
select 1 /*!00000AND 2>1*/
```

```
[('k', 'SELECT'), // keyword  
 ('1', '1'), // number  
 ('o', 'AND'), // operator  
 ('1', '2'), // number  
 ('o', '>'), // operator  
 ('1', '1')] // number
```

Meet the Tokens

- ▶ none/name
- ▶ variable
- ▶ string
- ▶ regular operator
- ▶ unknown
- ▶ number
- ▶ comment
- ▶ keyword
- ▶ group-like operation
- ▶ union-like operator
- ▶ logical operator
- ▶ function
- ▶ comma
- ▶ semi-colon
- ▶ left parens
- ▶ right parens

Merging, Specialization, Disambiguation

- ▶ "IS", "NOT" ==> "IS NOT" (single op)
- ▶ "NATURAL", "JOIN" => "NATURAL JOIN"
- ▶ "+" operator -> "+", "unary operator"
- ▶ COS, function, 1, number ==>
COS, not a function... not followed by (

Folding

- ▶ This step actually isn't needed to detect, but is needed to reduce false positives.
- ▶ Converts simple arithmetic expressions into a single value (don't try to evaluate them).
- ▶ 1-917-660-3400 -> "1"

Knows nothing about SQLi

- ▶ So far this is purely a parsing problem.
- ▶ Knows nothing about SQLi (which is evolving)
- ▶ Can be 100% tested against any SQL input (not SQLi) for correctness.

Fingerprints

- ▶ The token types of a user input form a hash or a fingerprint.
 - ▶ `select 1 /*!00000AND 2>1*/`
 - ▶ `K10101`
- ▶ Now let's generate fingerprints from Real World Data.
- ▶ Can we distinguish between SQLi and benign input?

Training on SQLi

- ▶ Parse known SQLi attacks from
 - ▶ SQLi vulnerability scanners
 - ▶ Published reports
 - ▶ SQLi How-Tos
- ▶ > 32,000 total

Training on real Input

- ▶ 100s of Millions of user inputs from Etsy's log were also parsed.
- ▶ Large enough to get a good sample (Top 50 USA site)
- ▶ Old enough to have lots of odd ways of handling query string, etc
- ▶ Full text search with an diverse subject domain

How many tokens are
needed to determine if a
user input
is SQLi or not?

5

No matter long the input is.

480 out of 1,048,576 are SQLi

```
n,(k( 1))Un n)ok( s)Unk 1)o1B n,(k1 s&n&s k(vv) sosos 1oks, sk)&1 1)o1f 1)o1k n);k& 1&1Bf 1)o1o s&os n,(kf s;k; n&1o1
1o(f( f(k,( 1ok1 loksc so1f( sk)&f soso( 1),(1 s))&( s))&1 &f()o nok(1 k1,1k 1&f(n soko1 1Unk1 1ok(1 n))kk 1Unk( soko(
1)o(1 s)ok1 sov&s n;kn( nok(k s))&f sovso 1)ok1 s))&o 1)ok( s&sos n&k(1 s&vso so1c sUk(k k1,1, 1)o(n 1)o(k 1Bf(1 1kf(1
s&ko1 s&k(o 1)k1 sk);k 1&f(1 1Unkf s)k1 sos&( 1&(k1 1))on s&kc nUk(o 1);ko( 1)B1 sokc n)o1& no1oo 1&(kn s&1oo so1o1
s)o1f 1&(kf s)o1k s)o1o f(1)o n&1o( s)o1B 1okv, sk)&( 1);kok sok1 f(k() 1Ukv, s&1of 1&1oo n&1f( 1)); 1))& 1&1of sovo1
s&1ov s)&1o sono1 1o((f 1))) 1o(s) s)&1f 1&1ov 1Uk1k n))ok k()ok nkksc 1Uk1c n))of s&(k1 s)&1B n;kks n)o(k kf(n, f(f(1
sovov s&1o( sovos s&1o1 vok1, sovok sUk1 1o((( 1)))k 1&1o1 f(f() 1)))o n))o( 1)))U k1k(k 1Uk1, 1&1f( so(s) 1)))B f(n(
n))o1 s)&(1 1)of( 1,(k1 sk)B1 f(1,f 1,(k 1Bk(1 lonos 1o1f( 1,f(1 1B1c s&okc s;ko( sk1os s&oko 1ono1 1,(kf sB1 1));k
s;kf( n)kks s;kok sk1o1 s;k(( 1o((1 1o1Bf so(f( n;kf( s&k(1 1&1o( nof(1 s);kk sk1c 1))o( s);ko s);kn sok1, s;k(1 1)kks
s);kf so(os so1ov s;k1, 1))Uk soknk s))k1 1)B1o 1)B1c n);k( n);kok s;k(o s);k( sok1o sok1c sf(n, s);k& sB1&s s;k1o sUno1
s))kk n);kf 1&so1 sokn, n;ko( n);kk n);kn n);ko s&1on sof(k n);k&k k1o(s sonos sk1&1 sof(f loso1 1;knc sUknk f()of n&(1
s&ko( sof() ok1o1 n,f(1 lo(1) s;kkn s;kks 1o(kn sof(1 sUkn, s)k1c 1;kn( s)k1o s;k&k skks s;n:k no(o1 s))o( k(ok( so(ks
so(kk so(kn so(ko s))o1 n)&(k o1kf( s))ok ;kknc skksc so(k1 n;k(( s&o(1 s))of so(k) n;k(1 n&(o1 s&kok sov:o s)of( sU(kk
sU(kn f(v,1 sk)of 1)&f( sk)ok no1f( sU(ks oUk1, 1ok1c s&(1) s&kos 1ok1k sUnk1 1)ono 1of(1 so1o( s;knn s;knk 1of() vUk1,
no1of 1&no1 sk)o1 s)B1 1)&o( sUk1& s&(k) 1o1)o f())&f sk)o( n&f(1 so1of 1)on& 1)B1& so1oo no1o1 so1ok 1ok1, 1of(n no1o(
so1os s;k( 1of(f sUnkf 1o(n) s&1os no(k1 n))o n))k 1kk(1 1;k(o 1)()s s&k1o s)B1& n)&1f n))&( sUk1, n)&1o no1&1 n));
sf(1) 1;k(1 n))& sokf( 1;k(( ook1, n)of( sUk1c s)B1c n&(k1 sUk1o s)B1o 1Ukf( okkkn s&vos s)o(k 1)&1f 1Uk1 1))&o 1))&f
1)&1B 1)&(k s,1), f(1o1 s)&f( s)o(1 sUkf( s&k&s 1okf( 1)&(1 1))&1 1;kf( 1))&( sokos 1))ok 1o1of 1o(1o 1kksc 1o1oo 1Uk(k
1))of 1o1ov Ukkkn 1,(f( 1ok(k so1Uk s&1f( sokok of(1) 1;k&k kf(1) sk)k1 s&v:o sok&s n)o1o n)o1f sUn(k 1o1o( 1o1o1 1))o1
sov&1 n));k n))&f sk)kk s)&(k 1)Unk n))&1 sU( (k 1)k1o 1);kk s;kvc 1);ko 1);kn 1)k1c s;kvk 1);kf 1Uks, s&o(k 1);k& s)&o(
s&(1o s&f() 1,1), 1);k( sk)Un sk)Uk s&f(1 1)&1o 1Uksc nUnk( so((k 1o1kf s&1Bf 1))kk kvk(1 n&o1o f(1)& &f(1 1))k1 so(((
s))Un s))Uk n,(f( 1)Uk1 s),(1 s&kkn 1))B1 s)kks 1Uk no(1) n)&f( s)ok( s))B1 sos 1&(1o s)Uk1 s));k so(1) 1&o(1 sok(1
nUk(k n&1of 1B1 sB1c n&1oo so(1o 1k1c sok(s sok(o sok(k so((s so1kf 1;kks s))B sf(s) 1&o1o n)k1o s))U sonk1 kf(1,
1o(kf 1,s), s))k so1&1 s))o s&nos s&1Uk s&o1o 1o(k1 so1Bf s;k[k sB1os of(o) s;k[n s))& s&(f( so1&s s&no1 so1&o s))));
```

Possible that more token types will be added to help reduce false positives.

The Library

- ▶ C, logic is under 1000 LOC
- ▶ No memory allocation
(caller makes a copy of input)
- ▶ Fixed, stable memory usage
- ▶ No threads
- ▶ 100k query strings can be checked per second
- ▶ Could go even faster

Sample Usage

```
sfilter sf; // on stack, ~500 bytes
const char* ucg = "my user input";
bool issqli = is_sql_i(&sf,
                      ucg, strlen(ucg));
// tada
```

metadata on input is in struct sfilter;

(names subject to change, cleanup)

Test Cases

- ▶ All input test cases available
- ▶ Including false positives found along the way
- ▶ Code coverage reports

Python Prototype

- ▶ Algorithm in python as well
- ▶ Not as up-to-date as the C version
- ▶ Working on it
- ▶ Runs under PyPy (and quite fast)

Make Existing Systems Work Better

- ▶ The Tokenizer could be ripped out, to make a "SQL normalizer/simplifier"
 - ▶ all white space normalized
 - ▶ all comments removed
 - ▶ all numbers in various flavors converted to "1"
 - ▶ all strings converted to a fixed value "foo"
- ▶ Makes existing regular expressions work better and detect more.

Great for Fuzzers

- ▶ The SQLi fingerprints are actually a great source of *templates* for fuzzers and SQLi generators
- ▶ Take fingerprint and turn it back into SQL

Available Now on GitHub

- ▶ <https://github.com/client9/libinjection>
- ▶ BSD License
(only to track how this gets used)
- ▶ Use it.

Help!

- ▶ More SQLi test cases!
- ▶ More real-world test cases
- ▶ Missing some PGSQL / Oracle string insanity
- ▶ Need better understanding of non-ASCII usage
- ▶ Porting to other languages (it's not that hard).

More Analysis at DEFCON 20

New Techniques in SQLi Obfuscation SQL never before used in SQLi

<http://slidesha.re/Mf0iNR>

July 27, 2012 Friday, 4:20pm at the Rio

Nick Galbreath



@ngalbreath



blackhat
USA 2012

<https://github.com/client9/libinjection>
<http://slidesha.re/0Bch5k>
@ngalbreath
nickg@client9.com
nickg@etsy.com

Thanks for coming by!