



Decision Procedures for Vulnerability Analysis

June 24, 2020

Benjamin Farinier

- Director Marie-Laure Potet
- Supervisor Sébastien Bardin
- Reviewer Thomas Jensen
- Reviewer Sylvain Conchon
- Examiner Mihaela Sighireanu
- Examiner Jean Goubault-Larrecq
- Examiner Roland Groz

Introduction Heartbleed

On April 7, 2014, the Heartbleed software vulnerability is made public

- between 24% and 55% of "secure" servers were affected
- due to a bug introduced on March 14, 2012 in the OpenSSL cryptography library
- allows an attacker to read secret data form the memory of a vulnerable server

The vulnerability is caused by a *buffer overflow*

- a kind of vulnerability known since 1972
- already exploited by the Morris computer worm in 1988

Long-standing causes known as dangerous are still the source of major vulnerabilities \implies need for research in formal verification



Formal verification aims to prove or disprove the correctness of a system with respect to a certain specification or property



Used in a growing number of contexts

- Cryptographic protocols
- Electronic hardware
- Software source code

Core concept: $\mathcal{M} \models \mathcal{P}$

- \mathcal{M} : the model of the system
- \mathcal{P} : the property to be checked
- = : the algorithmic check

Introduction Over-Approximation



correct behavioursincorrect behaviours

With an over-approximating model:

- All object behaviors are captured by the model
- But some model behaviors are not realizable by the object
- \Rightarrow False positive

Some over-approximation techniques:

Abstract Interpretation

object behaviours

model behaviours

- Hoare Logic
- ...



Introduction Under-Approximation



correct behavioursincorrect behaviours

With an under-approximating model:

- All model behaviors are realizable by the object
- But some object behaviors are not captured by the model
- $\Rightarrow \ \mathsf{False} \ \mathsf{negative}$

object behaviours model behaviours

Some under-approximation techniques:

- Bounded Model Checking (BMC)
- Symbolic Execution (SE)

• ...

Introduction Symbolic Execution



Symbolic Execution suffers several limitations...

- Path explosion
- Memory model
- Constraint solving
- Interactions with the environment

...but still leads to several successful applications

SAGE, P.Godefroid et al. \Rightarrow x86 instruction level SE $\label{eq:klee} \begin{array}{l} \mbox{KLEE, C.Cadar et al.} \\ \Rightarrow \mbox{LLVM bytecode level SE} \end{array}$

It is now a question of applying it to vulnerability analysis

Introduction Motivating Example

```
#define SIZE
```

```
void get_secret (char secr[]) {
// Retrieve the secret
}
```

```
void read_input (char src[], char dst[]) {
    int i = 0;
    while (src[i]) {
        dst[i] = src[i];
        i++;
    }
}
```

```
int validate (char secr[], char inpt[]) {
    int b = 1;
    for (int i = 0; i < SIZE; i++) {
        b &= secr[i] == inpt[i];
    }
    return b;
}</pre>
```

```
int main (int argc, char *argv[]) {
    char secr[SIZE];
    char inpt[SIZE];
    if (argc != 2) return 0;
    get_secret(secr);
    read_input(argv[1],inpt);
    if (validate(secr, inpt)) {
        printf("Success!\n");
        }
    else {
        printf("Failure...\n");
        }
    }
}
```

Introduction Motivating Example

```
#define SIZE
```

```
void get_secret (char secr[]) {
// Retrieve the secret
}
```

```
void read_input (char src[], char dst[]) {
    int i = 0;
    while (src[i]) {
        dst[i] = src[i];
        i++;
    }
}
```

```
int validate (char secr[], char inpt[]) {
    int b = 1;
    for (int i = 0; i < SIZE; i++) {
        b &= secr[i] == inpt[i];
    }
    return b;
}</pre>
```

```
int main (int argc, char *argv[]) {
    char secr[SIZE];
    char inpt[SIZE];
    if (argc != 2) return 0;
    get_secret(secr);
    read_input(argv[i],inpt);
    if (validate(secr, inpt)) {
        printf("Success!\n");
      }
    else {
        printf("Failure...\n");
    }
}
```

Goal

Find an input such that the execution reach the "Success!" branch

Introduction Motivating Example

```
#define SIZE
```

}

```
void get_secret (char secr[]) {
// Retrieve the secret
}
```

```
void read_input (char src[], char dst[]) {
    int i = 0;
    while (src[i]) {
        dst[i] = src[i];
        i++;
    }
}
```

```
int validate (char secr[], char inpt[]) {
    int b = 1;
    for (int i = 0; i < SIZE; i++) {
        b &= secr[i] == inpt[i];
    }
    return b;
}</pre>
```

```
int main (int argc, char *argv[]) {
    char secr[SIZE];
    char inpt[SIZE];
    if (argc != 2) return 0;
    get_secret(secr);
    read_input(argv[1],inpt);
    if (validate(secr, inpt)) {
        printf("Success!\n");
      }
    else {
        printf("Failure...\n");
    }
}
```

```
 \exists i. \exists s. \exists m_0. \exists p_0. \qquad i: \text{ input } m: \text{ memory} \\ p_1 \triangleq p_0 - \text{SIZE} \qquad s: \text{ secret } p: \text{ stack pointer} \\ p_2 \triangleq p_1 - \text{SIZE} \\ m_1 \triangleq m_0 [p_1 \dots p_1 + \text{SIZE} - 1] \leftarrow s \\ m_2 \triangleq m_1 [p_2 \dots p_2 + N - 1] \leftarrow i \\ m_2 [p_1 \dots p_1 + \text{SIZE} - 1] = m_2 [p_2 \dots p_2 + \text{SIZE} - 1]
```





Unrolling-based verification techniques (BMC, SE)

- may produce huge formulas
- with a high number of reads and writes

In some extreme cases, solvers may spend hours on these formulas



ASPack case study: 293 000 reads, 58 000 writes \Rightarrow 24 hours of resolution !

Sending the formula to a solver:

$$\Rightarrow \ \left\{ \textit{s}_{[0 \ . \ \text{size}-1]} = 0, \textit{i}_{[0 \ . \ \text{size}-1]} = 0, \ldots \right\}$$

"If the secret is 0, then you can choose 0 as an input."

Sure, that is true... but a false positive in practice

- the secret will not likely be 0
- \Rightarrow the execution will not reach the "Success" branch

Sending the formula to a solver:

$$\Rightarrow \ \left\{ \textit{s}_{[0 \ . \ \text{size}-1]} = 0, \textit{i}_{[0 \ . \ \text{size}-1]} = 0, \ldots \right\}$$

"If the secret is 0, then you can choose 0 as an input."

Sure, that is true... but a false positive in practice

- the secret will not likely be 0
- \Rightarrow the execution will not reach the "Success" branch

Threat models make security \neq safety

A better formalization:

- We do not have control over s, m_0 and p_0
- These variables should be universally quantified
- \Rightarrow This is where the problems begin...

- Symbolic Execution (SE)
 - $\circ~$ under-approximation verification technique
 - heavily relies on SMT solvers
- Application to vulnerability analysis
 - $\circ\;$ requires to move from source analysis to binary analysis
 - $\circ~$ modeling threat models introduces universal quantifiers
- Problems
 - finding a model for a \forall -formula is difficult
 - o going low-level significantly increases formula size
 - \Rightarrow The Death of SMT Solvers

O Introduction

- 1 Model Generation for Quantified Formulas: A Taint-Based Approach
- 2 Arrays Made Simpler: An Efficient, Scalable and Thorough Preprocessing
- 3 Get Rid of False Positives with Robust Symbolic Execution

4 Conclusion

Section 1

Model Generation for Quantified Formulas: A Taint-Based Approach

- Challenge
 - Deal with quantified-formulas and model generation
 - Notoriously hard! (undecidable)
- Existing approaches
 - Complete but costly for very specific theories
 - Incomplete but efficient for UNSAT/UNKNOWN
 - Costly or too restricted for model generation
- Our proposal
 - $\circ~{\rm SAT}/{\rm UNKNOWN}$ and model generation
 - $\circ~$ Incomplete but efficient, generic, theory independent
 - $\circ~$ Reuse state-of-the-art solvers as much as possible

Published in Computer Aided Verification 30th, Oxford, UK, 2018 [CAV18] Presented in Approches Formelles dans l'Assistance au Développement de Logiciels, Grenoble, France, 2018 [AFADL18]

```
int main () {
    int a = input ();
    int b = input ();
    int x = rand ();
    if (a * x + b > 0) {
        analyze_me();
    }
    else {
        ...
    }
}
```

We propose a way to infer such conditions

- Quantified reachability condition:
 ∀x.ax + b > 0
- Generalizable solutions of ax + b > 0 have to be independent from x
 - A bad solution:

$$a=1\wedge x=1\wedge b=0$$

• A good solution:
$$a = 0 \land x = 1 \land b = 1$$

- The constraint *a* = 0 is the independence condition
- Quantifier-free reachability condition: (ax + b > 0) ∧ (a = 0)

Model Generation for Quantified Formulas Our Proposal in a Nutshell



Sufficient Independence Condition (SIC)

A SIC for a formula $\Phi(\mathbf{x}, \mathbf{a})$ with regard to \mathbf{x} is a formula $\Psi(\mathbf{a})$ such that $\Psi(\mathbf{a}) \models (\forall \mathbf{x}.\forall \mathbf{y}.\Phi(\mathbf{x}, \mathbf{a}) \Leftrightarrow \Phi(\mathbf{y}, \mathbf{a})).$

— formula indep.

- If $\Phi \triangleq ax + b > 0$ then a = 0 is a $SIC_{\Phi,x}$.
- If $\Delta \triangleq (t [a] \leftarrow b) [c]$ then a = c is a $SIC_{\Delta,t}$.
- \perp is always a SIC, but a useless one...

Model generalization

- Let $\Phi(\mathbf{x}, \mathbf{a})$ a formula and $\Psi(\mathbf{a})$ a $SIC_{\Phi, \mathbf{x}}$.
- If there exists an interpretation {x, a} such that {x, a} ⊨ Ψ (a) ∧ Φ (x, a), then {a} ⊨ ∀x.Φ (x, a).

Weakest Independence Condition (WIC)

A WIC for a formula $\Phi(x, a)$ with regard to x is a SIC_{Φ,x} Π such that, for any other SIC_{Φ,x} $\Psi, \Psi \models \Pi$.

- Both SIC a = 0 and a = c presented earlier are WIC.
- $\Omega \triangleq \forall x. \forall y. (\Phi(x, a) \Leftrightarrow \Phi(y, a))$ is always a $WIC_{\Phi,x}$, but involves quantifiers
- A formula Π is a $\operatorname{WIC}_{\Phi, \mathbf{x}}$ if and only if $\Pi \equiv \Omega$.

Model specialization

- Let $\Phi(\mathbf{x}, \mathbf{a})$ a formula and $\Pi(\mathbf{a})$ a $\operatorname{WIC}_{\Phi, \mathbf{x}}$.
- If there exists an interp. {a} such that {a} $\models \forall x.\Phi(x, a)$, then $\{x, a\} \models \Pi(a) \land \Phi(x, a)$ for any valuation x of x.

```
Function inferSIC(\phi.x):
     Input: \Phi a formula and x a set of targeted variables
     Output: \Psi a SIC<sub>\Phi,x</sub>
     either \Phi is a constant
      l return ⊤
     either \Phi is a variable v
                                                                                          syntactic part
      _ return v ∉ x
                                                                              a and b indep<sub>x</sub> \rightsquigarrow f (a, b) indep<sub>x</sub>
     either \Phi is a function f(\phi_1, .., \phi_n)
           Let \psi_i \triangleq \texttt{inferSIC}(\phi_i, \mathbf{x}) for all i \in \{1, ..., n\}
           Let \Psi \triangleq \text{theorySIC}(f, (\phi_1, .., \phi_n), (\psi_1, .., \psi_n), \mathbf{x})
           return \Psi \vee \bigwedge_{i} \psi_{i}
                                                                                          semantic part
                                                                             a indep<sub>x</sub> and a = 0 \rightsquigarrow a \cdot * indep_x
```

Proposition

 If theorySIC(f, φ_i, ψ_i, x) computes a SIC_f(φ_i),x, then inferSIC(Φ, x) computes a SIC_{Φ,x}.

```
      Function inferSIC((\Phi, x)):

      Input: \Phi a formula and x a set of targeted variables

      Output: \Psi a SIC_{\Phi,x}

      either \Phi is a constant

      \Box return \top

      either \Phi is a variable v

      syntactic part

      \Box return v \notin x

      either \Phi is a function f(\phi_1, .., \phi_n)

      Let \psi_i \triangleq inferSIC(\phi_i, x) for all i \in \{1, .., n\}

      Let \Psi \triangleq theorySIC(f, (\phi_1, .., \phi_n), (\psi_1, .., \psi_n), x)

      return \Psi \lor \bigwedge_i \psi_i

      semantic part

      a indep_x and a = 0 \rightsquigarrow a \cdot * indep_x
```

theorySIC defined as a recursive function

$$(a \Rightarrow b)^{\bullet} \triangleq (a^{\bullet} \land a = \bot) \lor (b^{\bullet} \land b = \top)$$
$$(a \land b)^{\bullet} \triangleq (a^{\bullet} \land a = \bot) \lor (b^{\bullet} \land b = \bot)$$
$$(a \lor b)^{\bullet} \triangleq (a^{\bullet} \land a = \top) \lor (b^{\bullet} \land b = \top)$$
$$(\text{ite } c a b)^{\bullet} \triangleq (c^{\bullet} \land \text{ite } c a^{\bullet} b^{\bullet}) \lor (a^{\bullet} \land b^{\bullet} \land a = b)$$

$$(a_n \wedge b_n)^{\bullet} \triangleq (a_n^{\bullet} \wedge a_n = 0_n) \vee (b_n^{\bullet} \wedge b_n = 0_n) (a_n \vee b_n)^{\bullet} \triangleq (a_n^{\bullet} \wedge a_n = 1_n) \vee (b_n^{\bullet} \wedge b_n = 1_n) (a_n \times b_n)^{\bullet} \triangleq (a_n^{\bullet} \wedge a_n = 0_n) \vee (b_n^{\bullet} \wedge b_n = 0_n) (a_n \ll b_n)^{\bullet} \triangleq (b_n^{\bullet} \wedge b_n \ge n)$$

$$((a[i] \leftarrow e)[j])^{\bullet} \triangleq (\text{ite } (i = j) \ e \ (a[j]))^{\bullet} \\ \triangleq ((i = j)^{\bullet} \land (\text{ite } (i = j) \ e^{\bullet} \ (a[j])^{\bullet})) \\ \lor (e^{\bullet} \land (a[j])^{\bullet} \land (e = a[j])) \\ \triangleq (i^{\bullet} \land j^{\bullet} \land (\text{ite } (i = j) \ e^{\bullet} \ (a[j])^{\bullet})) \\ \lor (e^{\bullet} \land (a[j])^{\bullet} \land (e = a[j]))$$

Model Generation for Quantified Formulas Experimental Evaluation

Boolector: an	efficient	$\operatorname{QF}\operatorname{-solver}$	for
 bitvectors and	arrays		

Best	approaches	
------	------------	--

			J.	
		Z3	Btor	Btor● ⊳ Z3
IB	SAT	261	399	485
Ę	# UNSAT	165	N/A	165
Ę	UNKNOWN	843	870	619
\mathbf{S}	total time	270 150	350	94 610
0	SAT	953	1042	1067
SE	# UNSAT	319	N/A	319
BIN	UNKNOWN	149	379	35
	total time	64761	1 1 5 2	1 169

GRUB example

	Z3	Btor●
SAT	1	540
# UNSAT	42	N/A
UNKNOWN	852	355
total time	159 765	16732

Complementarity with existing solvers (SAT instances)

		CVC4•	Z3•	Btor●
SMT-LIB	CVC4	-10 +168 [252]		-10 +325 [409]
	Z3		-119 +224 [485]	-86 +224 [485]
BINSEC	CVC4	-25 +28 [979]		-25 +116 [1067]
DINSEC	Z3		-25 +114 [1067]	-25 +114 [1067]

solver•: solver enhanced with our method

Section 2

Arrays Made Simpler: An Efficient, Scalable and Thorough Preprocessing

Challenge

 Array theory useful for modelling memory or data structures... code ...but a bottleneck for resolution of large formulas (BMC, SE) Existing approaches • General decision procedures for the theory of arrays • Dedicated handling of arrays inside tools formula Our proposal SMT • FAS, an efficient simplification for array theory \Rightarrow Improves existing solvers solver Published in Logic for Programming, Artificial Intelligence and Reasoning, Awassa, Ethiopia, 2018 [LPAR18] Presented in Journées Francophones des Langages Applicatifs, Y/N Banyuls-sur-Mer, France, 2018 [JFLA18]

Two basic operations on arrays

- Reading in *a* at index $i \in \mathcal{I}$: *a*[*i*]
- Writing in a an element $e \in \mathcal{E}$ at index $i \in \mathcal{I}$: $a[i] \leftarrow e$

$$\begin{array}{c} \cdot [\cdot] : \operatorname{Array} \ \mathcal{I} \ \mathcal{E} \to \mathcal{I} \to \mathcal{E} \\ \cdot [\cdot] \leftarrow \cdot : \operatorname{Array} \ \mathcal{I} \ \mathcal{E} \to \mathcal{I} \to \mathcal{E} \to \operatorname{Array} \ \mathcal{I} \ \mathcal{E} \end{array}$$

$$\operatorname{ROW-axiom:} \ \forall a \ i \ j \ e. \ (a \ [i] \leftarrow e) \ [j] = \begin{cases} e & \text{if } i = j \\ a \ [j] & \text{otherwise} \end{cases}$$

Prevalent in software analysis

- Modelling memory
- Abstracting data structure (map, queue, stack...)

Hard to solve

- NP-complete
- Read-Over-Write (ROW) may require case-splits

Unrolling-based verification techniques (BMC, SE)

- may produce huge formula
- high number of reads and writes

In some extremes cases, solvers may spend hours on these formulas

Without proper simplification, array theory might become a bottleneck for resolution

What should we simplify ? Read-Over-Write (ROW)!

An example coming from binary analysis

```
esp<sub>0</sub> : BitVec16
mem<sub>0</sub> : Array BitVec16 BitVec16
```

```
\begin{array}{l} \texttt{assert} (\texttt{esp}_0 > 61440) \\ \texttt{mem}_1 \triangleq \texttt{mem}_0 [\texttt{esp}_0 - 16] \leftarrow 1415 \\ \texttt{esp}_1 \triangleq \texttt{esp}_0 - 64 \\ \texttt{eax}_0 \triangleq \texttt{mem}_1 [\texttt{esp}_1 + 48] \\ \texttt{assert} (\texttt{mem}_1 [\texttt{eax}_0] = 9265) \end{array}
```

esp₀ : BitVec16
mem₀ : Array BitVec16 BitVec16

$assert(esp_0$	> 61440)
assert (mem ₀	[1415] = 9265)

These simplifications depend on two factors

The equality check procedure

verify that $esp_1 + 48 = esp_0 - 16$

- \Rightarrow precise reasoning: base normalization + abstract domains
- The underlying representation of an array

remember that $mem_1 [esp_1 + 48] = 1415$

⇒ scalability issue: list-map representation

Arrays Made Simpler Improving scalability: list-map representation



How to update

Given a write of e at index i

- Is *i comparable* with indices of elements in the head?
- If so add (*i*, *e*) in this map
- Else add a new head map containing only (*i*, *e*)

How to simplify ROW

Given a read at index j

- Is *j* comparable with indices of elements in the head?
- If so, look for (*i*, *e*) with *i*=*j*
 - $\circ~$ if succeeds then return e
 - else recurse on next map
- Else stop

Propagate "variable+constant" terms

- If $y \triangleq z+1$ then $x \triangleq y+2 \rightsquigarrow x \triangleq z+3$
- Together with associativity, commutativity...
- $\Rightarrow\,$ Reduce the number of bases

Associate to every indices i an abstract domain i^{\sharp}

- If $i^{\sharp} \sqcap j^{\sharp} = \bot$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation
- \Rightarrow Prove disequality between different bases

- + 6,590 \times 3 medium-size formulas from static SE
- TIMEOUT = 1,000 seconds

		simpl.		#TIME	#POW				
		time	Bo	Boolector Yices Z3				#ROW	
ete	default	61	0	163	2	69	0	872	866,155
JCre	FAS	85	0	94	2	68	0	244	1,318
cor	FAS-itv	111	0	94	2	68	0	224	1,318
interval	default	65	0	2,584	2	465	31	155,992	866,155
	FAS	99	0	2,245	2	487	25	126,806	531,654
	FAS-itv	118	0	755	2	140	14	37,269	205,733
olic	default	61	0	6,173	3	1,961	65	305,619	866,155
q	FAS	91	0	6,117	3	1,965	66	158,635	531,654
syn	FAS-itv	111	0	4,767	2	1,108	43	80,569	295,333

- 29 \times 3 very large formulas from dynamic SE
- TIMEOUT = 1,000 seconds

simpl. #TIMEOUT and resolution time						#POW				
		time	Boo	olector	Ň	ices	Z3		#-100	
e	default	44	10	159	4	1,098	26	3.33	1,120,798	
cret	FAS-list	1,108	8	845	4	198	10	918	456,915	
ouo	FAS	196	8	820	4	196	10	922	456,915	
Ŭ	FAS-itv	210	4	654	1	12	4	1,120	0	
_	default	44	12	131	12	596	27	0.19	1,120,798	
Z	FAS-list	222	12	129	12	595	26	236	657,594	
nte	FAS	231	12	129	12	597	26	291	657,594	
	FAS-itv	237	12	58	12	28	19	81	651,449	
. <u>∪</u>	default	40	12	1,522	12	1,961	27	0.13	1,120,798	
0	FAS-list	187	11	1,199	12	2,018	26	486	657,594	
Т,	FAS	194	11	1,212	12	2,081	26	481	657,594	
ίν'	FAS-itv	200	11	1,205	12	2,063	26	416	657,594	

Arrays Made Simpler Focus on Specific Case: the ASPack Example

- Huge formula obtained from the ASPack packing tool
- 293000 ROWs
- 24 hours of resolution!



Using FAS

- #ROW reduced to 2467
- 14 sec for resolution
- 61 sec for preprocessing

Using list representation

- Same result with a bound of 385 024 and beyond...
- ...but 53 min preprocessing

Section 3

Get Rid of False Positives with Robust Symbolic Execution

- Symbolic Execution (SE)
 - under-approximation verification technique
 - heavily relies on SMT solvers
 - $\circ~$ should be exempt of false positives
- In practice, false positives exist
 - misspecified abstractions, initial state...
 - some ad hoc workarounds, no real solution
- Our proposal: Robust Symbolic Execution
 - o distinguish between controlled and uncontrolled inputs
 - $\circ\;$ robust solutions are independent of uncontrolled inputs
 - practical application of [CAV18] and [LPAR18]

Presented in Journées Francophones des Langages Applicatifs, Les Rousses, France, 2019 [JFLA19]

Robust Symbolic Execution Motivating Example Remembered

```
#define SIZE
void get_secret (char secr[]) {
// Retrieve the secret
3
                                                        int main (int argc, char *argv[]) {
                                                           char secr[SIZE]:
void read_input (char src[], char dst[]) {
                                                           char inpt[SIZE];
  int i = 0;
  while (src[i]) {
                                                           if (argc != 2) return 0:
    dst[i] = src[i]:
    i++;
                                                           get_secret(secr);
 }
                                                           read_input(argv[1], inpt);
3
                                                           if (validate(secr, inpt)) {
int validate (char secr[], char inpt[]) {
                                                             printf("Success!\n");
  int b = 1:
                                                           3
  for (int i = 0; i < SIZE; i++) {</pre>
                                                           else {
     b &= secr[i] == inpt[i];
                                                             printf("Failure...\n");
                                                           3
  }
                                                         }
  return b:
3
  \exists i. \exists s. \exists m_0. \exists p_0.
                                                 i: input m: memory
      p_1 \triangleq p_0 - \text{SIZE}
                                                 s: secret p: stack pointer
      p_2 \triangleq p_1 - \text{SIZE}
     m_1 \triangleq m_0 [p_1 \dots p_1 + \text{SIZE} - 1] \leftarrow s
     m_2 \triangleq m_1 [p_2 \dots p_2 + N - 1] \leftarrow i
  m_2[p_1 \dots p_1 + \text{SIZE} - 1] = m_2[p_2 \dots p_2 + \text{SIZE} - 1]
```

June 24, 2020 — Benjamin Farinier — 33/41

$$\begin{array}{l} \exists i.\exists s.\exists m_0.\exists p_0.\\ p_1 \triangleq p_0 - \text{SIZE}\\ p_2 \triangleq p_1 - \text{SIZE}\\ m_1 \triangleq m_0 \left[p_1 \dots p_1 + \text{SIZE} - 1 \right] \leftarrow s\\ m_2 \triangleq m_1 \left[p_2 \dots p_2 + N - 1 \right] \leftarrow i\\ m_2 \left[p_1 \dots p_1 + \text{SIZE} - 1 \right] = m_2 \left[p_2 \dots p_2 + \text{SIZE} - 1 \right] \end{array}$$

Sending the formula to a solver:

$$\Rightarrow \left\{ s_{[0 \dots \text{SIZE}-1]} = 0, i_{[0 \dots \text{SIZE}-1]} = 0, \dots \right\}$$

• This is a false positive

A better formalization: Robust SE

- We do not have control over s, m_0 and p_0
- These variables should be universally quantified

$$\exists i.\forall s.\forall m_0.\forall p_0.$$

$$p_1 \triangleq p_0 - \text{SIZE}$$

$$p_2 \triangleq p_1 - \text{SIZE}$$

$$m_1 \triangleq m_0 [p_1 \dots p_1 + \text{SIZE} - 1] \leftarrow s$$

$$m_2 \triangleq m_1 [p_2 \dots p_2 + N - 1] \leftarrow i$$

$$m_2 [p_1 \dots p_1 + \text{SIZE} - 1] = m_2 [p_2 \dots p_2 + \text{SIZE} - 1]$$

Problems:

- finding a model for a ∀-formula is difficult
- going low-level significantly increases formula size
- $\Rightarrow\,$ The Death of SMT Solvers

```
\begin{aligned} \exists i. \exists s. \exists m_0. \forall p_0. \\ p_1 &\triangleq p_0 - \text{SIZE} \\ p_2 &\triangleq p_1 - \text{SIZE} \\ m_1 &\triangleq m_0 [p_0 - \text{SIZE} \dots p_0 - 1] \leftarrow s \\ m_2 &\triangleq m_1 [p_0 - 2 \cdot \text{SIZE} \dots p_0 - 2 \cdot \text{SIZE} + N - 1] \leftarrow i \\ i[0 \dots \text{SIZE} - 1] &= i[\text{SIZE} \dots 2 \cdot \text{SIZE} - 1] \\ \land N \geq 2 \cdot \text{SIZE} \end{aligned}
```

Problems:

- finding a model for a ∀-formula is difficult [CAV18]
- going low-level significantly increases formula size [LPAR18]
- \Rightarrow The Death of SMT Solvers

For example with SIZE = 8,

- input abcdefghabcdefgh leads to the "Success!" branch
- buffer overflow in read_input

• cot of	ara alima ch	allongos		SE robust			
• Set Of		lalleliges	true	false			
 comparison 	re true and	false positi	positives	positives	UNKNOWN		
				Boolector	N/A	N/A	24
				CVC4	5	0	19
				Yices	N/A	N/A	24
				Z3	7	0	17
							1.
	SE classic				SE	robust + e	elim.
	true	false			true	false	
	positives	positives	UNKNOWN		positives	positives	UNKNOWN
Boolector	12	11	1	Boolector	12	0	12
CVC4	7	9	8	CVC4	7	0	17
Yices	7	11	6	Yices	7	0	17
70	10	10	0	73	12	0	12

Back to 28: GRUB2 Authentication Bypass

- Original version: press Backspace 28 times to get a rescue shell
- Case study: same vulnerable code turned into a crackme challenge
- SE classic: incorrect solution
- SE robust: solvers TIMEOUT

- SE robust + elim.: correct solution in 80s
- SE robust + elim. + simpl.: correct solution in 30s

Section 4

Conclusion

- Symbolic Execution (SE)
 - $\circ~$ under-approximation verification technique
 - heavily relies on SMT solvers
- Application to vulnerability analysis
 - $\circ\;$ requires to move from source analysis to binary analysis
 - $\circ~$ modeling threat models introduces universal quantifiers
- Problems
 - finding a model for a \forall -formula is difficult
 - o going low-level significantly increases formula size
 - \Rightarrow The Death of SMT Solvers

Model Generation for Quantified Formulas

- · Proposed a novel and generic taint-based approach
- Proved its correctness and its efficiency
- $\circ~$ Presented an implementation for arrays and bit-vectors
- $\circ~$ Evaluated on $\rm SMT\text{-}LIB$ and formulas generated by Symbolic Execution

Arrays Made Simpler

- Presented FAS, a simplification dedicated to the theory of arrays
- Geared at eliminating ROW, based on a dedicated data structure, original simplifications and low-cost reasoning
- Evaluated in different settings on very large formulas

8 Robust Symbolic Execution

- Highlighted the problem of false positives in classic Symbolic Execution
- Introduced formally the framework of Robust Symbolic Execution
- $\circ~$ Implemented a proof of concept in the binary analyser $\rm BINSEC$

Model Generation for Quantified Formulas

- o More precise inference mechanisms of independence conditions
- $\circ~$ Identification of subclasses for which inferring WIC is feasible
- o Combination with other quantifier instantiation techniques

Arrays Made Simpler

- Deeper integration inside a dedicated array solver
- Adding more expressive domain reasoning

8 Robust Symbolic Execution

- $\circ~$ Precise evaluation of our semi-automatic incremental specification procedure
- $\circ~$ Thorough comparison of Robust Symbolic Execution to other techniques

Beyond that

- $\circ~$ Restrict to "there exists" and "for all" quantifications is not nuanced enough
- $\circ~$ Might want to say that an event "almost always" or "almost never" occurs
- $\Rightarrow\,$ Requires notions of probabilities or model counting

Conclusion Bibliography



Benjamin Farinier, S. Bardin, R. Bonichon, and M. Potet.

Model generation for quantified formulas: A taint-based approach.

In Computer Aided Verification - 30th International Conference, CAV 2018, Oxford, UK, July 14-17, 2018, 2018.



Benjamin Farinier, R. David, S. Bardin, and M. Lemerre.

Arrays made simpler: An efficient, scalable and thorough preprocessing. In LPAR-22. 22nd International Conference on Logic for Programming, Artificial Intelligence and Reasoning, Awassa, Ethiopia, November 16-21, 2018, 2018.



Benjamin Farinier, S. Bardin, R. Bonichon, and M. Potet.

Génération de modèles pour les formules quantifiées : une approche basée sur la teinte. In Approches Formelles dans l'Assistance au Développement de Logiciels, AFADL 2018, Grenoble, France, June 13-15, 2018, 2018.



Benjamin Farinier, R. David, and S. Bardin.

Simplification efficace pour la théorie des tableaux.

In Journées Francophones des Langages Applicatifs, JFLA 2018, Banyuls-sur-Mer, France, January 24-27, 2018, 2018.



Benjamin Farinier, S. Bardin, R. Bonichon, and M. Potet.

En finir avec les faux positifs grâce à l'exécution symbolique robuste.

In Journées Francophones des Langages Applicatifs, JFLA 2019, Les Rousses, France, January 30-February 2, 2019, 2019.