



Decision Procedures for Vulnerability Analysis

Journées du GDR GPL 2021

Benjamin Farinier

Director — Marie-Laure Potet
Supervisor — Sébastien Bardin

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
- \models : the algorithmic check

Some automated software verification techniques:

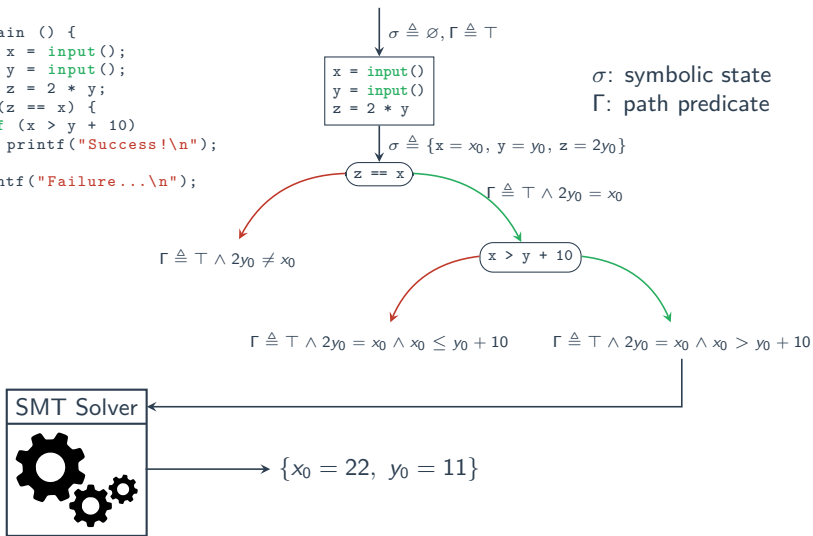
- Abstract Interpretation
- Bounded Model Checking (BMC)
- **Symbolic Execution (SE)**



Introduction

Symbolic Execution

```
int main () {  
  int x = input();  
  int y = input();  
  int z = 2 * y;  
  if (z == x) {  
    if (x > y + 10)  
      printf("Success!\n");  
  }  
  printf("Failure...\n");  
}
```



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.

⇒ x86 instruction level SE



KLEE, C.Cadar et al.

⇒ LLVM bytecode level SE

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;
}

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;
}

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");
    }
}
```

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;
}

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");
    }
}
```

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

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

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

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

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

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

i : input m : memory

s : secret p : stack pointer

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

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

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

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

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

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

oversimplified formula!

The real formula is about 2130 reads and 456 writes

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

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

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

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

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

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

oversimplified formula!

The real formula is about 2130 reads and 456 writes

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

⇒ 24 hours of resolution !

Sending the formula to a solver:

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

“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:

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

“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)
 - under-approximation verification technique
 - heavily relies on SMT solvers
- Application to vulnerability analysis
 - requires to move from source analysis to binary analysis
 - modeling threat models introduces universal quantifiers
- Problems
 - finding a model for a \forall -formula is difficult
 - going low-level significantly increases formula size
 - ⇒ The Death of SMT Solvers

- 0 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
 - SAT/UNKNOWN and model generation
 - Incomplete but efficient, generic, theory independent
 - 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](#)]

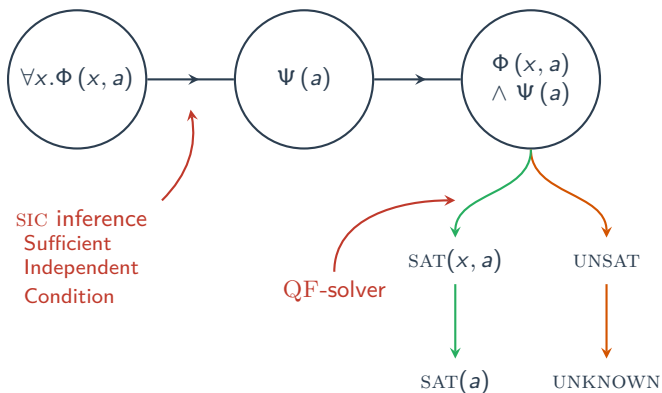
```
int main () {  
    int a = input ();  
    int b = input ();  
  
    int x = rand ();  
  
    if (a * x + b > 0) {  
        analyze_me ();  
    }  
    else {  
        ...  
    }  
}
```

We propose a way to infer
such conditions

- Quantified reachability condition:
 $\forall x. ax + b > 0$
- 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:
 $a = 0 \wedge x = 1 \wedge b = 1$
- The constraint $a = 0$ is the **independence condition**
- Quantifier-free reachability condition:
 $(ax + b > 0) \wedge (a = 0)$

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 $\Psi(\mathbf{a}) \models (\forall \mathbf{x}. \forall \mathbf{y}. \Phi(\mathbf{x}, \mathbf{a}) \Leftrightarrow \Phi(\mathbf{y}, \mathbf{a}))$.

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

formula indep.

Model generalization

- Let $\Phi(\mathbf{x}, \mathbf{a})$ a formula and $\Psi(\mathbf{a})$ a $\text{SIC}_{\Phi, \mathbf{x}}$.
- If there exists an interpretation $\{\mathbf{x}, \mathbf{a}\}$ such that $\{\mathbf{x}, \mathbf{a}\} \models \Psi(\mathbf{a}) \wedge \Phi(\mathbf{x}, \mathbf{a})$, then $\{\mathbf{a}\} \models \forall \mathbf{x}. \Phi(\mathbf{x}, \mathbf{a})$.

Weakest Independence Condition (WIC)

A WIC for a formula $\Phi(\mathbf{x}, \mathbf{a})$ with regard to \mathbf{x} is a $\text{SIC}_{\Phi, \mathbf{x}}$ Π such that, for any other $\text{SIC}_{\Phi, \mathbf{x}}$ Ψ , $\Psi \models \Pi$.

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

Model specialization

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

Function $\text{inferSIC}(\Phi, \mathbf{x})$:

Input: Φ a formula and \mathbf{x} a set of targeted variables

Output: Ψ a $\text{SIC}_{\Phi, \mathbf{x}}$

either Φ is a constant

└ **return** \top

either Φ is a variable v

└ **return** $v \notin \mathbf{x}$

either Φ is a function $f(\phi_1, \dots, \phi_n)$

└ Let $\psi_i \triangleq \text{inferSIC}(\phi_i, \mathbf{x})$ for all $i \in \{1, \dots, n\}$

└ Let $\Psi \triangleq \text{theorySIC}(f, (\phi_1, \dots, \phi_n), (\psi_1, \dots, \psi_n), \mathbf{x})$

└ **return** $\Psi \vee \bigwedge_i \psi_i$

syntactic part

a and $b \text{ indep}_x \rightsquigarrow f(a, b) \text{ indep}_x$

semantic part

$a \text{ indep}_x$ and $a = 0 \rightsquigarrow a \cdot * \text{ indep}_x$

Proposition

- If $\text{theorySIC}(f, \phi_i, \psi_i, \mathbf{x})$ computes a $\text{SIC}_{f(\phi_i), \mathbf{x}}$, then $\text{inferSIC}(\Phi, \mathbf{x})$ computes a $\text{SIC}_{\Phi, \mathbf{x}}$.

Function $\text{inferSIC}(\Phi, \mathbf{x})$:

Input: Φ a formula and \mathbf{x} a set of targeted variables

Output: Ψ a $\text{SIC}_{\Phi, \mathbf{x}}$

either Φ is a constant

└ **return** \top

either Φ is a variable v

└ **return** $v \notin \mathbf{x}$

either Φ is a function $f(\phi_1, \dots, \phi_n)$

└ Let $\psi_i \triangleq \text{inferSIC}(\phi_i, \mathbf{x})$ for all $i \in \{1, \dots, n\}$

└ Let $\Psi \triangleq \text{theorySIC}(f, (\phi_1, \dots, \phi_n), (\psi_1, \dots, \psi_n), \mathbf{x})$

└ **return** $\Psi \vee \bigwedge_i \psi_i$

syntactic part

a and $b \text{ indep}_x \rightsquigarrow f(a, b) \text{ indep}_x$

semantic part

$a \text{ indep}_x$ and $a = 0 \rightsquigarrow a \cdot * \text{ indep}_x$

theorySIC defined as a recursive function

$$\begin{aligned}(a \Rightarrow b)^\bullet &\triangleq (a^\bullet \wedge a = \perp) \vee (b^\bullet \wedge b = \top) \\(a \wedge b)^\bullet &\triangleq (a^\bullet \wedge a = \perp) \vee (b^\bullet \wedge b = \perp) \\(a \vee b)^\bullet &\triangleq (a^\bullet \wedge a = \top) \vee (b^\bullet \wedge b = \top) \\(\text{ite } c \ a \ b)^\bullet &\triangleq (c^\bullet \wedge \text{ite } c \ a^\bullet \ b^\bullet) \vee (a^\bullet \wedge b^\bullet \wedge a = b)\end{aligned}$$

$$\begin{aligned}(a_n \wedge b_n)^\bullet &\triangleq (a_n^\bullet \wedge a_n = 0_n) \vee (b_n^\bullet \wedge b_n = 0_n) \\(a_n \vee b_n)^\bullet &\triangleq (a_n^\bullet \wedge a_n = 1_n) \vee (b_n^\bullet \wedge b_n = 1_n) \\(a_n \times b_n)^\bullet &\triangleq (a_n^\bullet \wedge a_n = 0_n) \vee (b_n^\bullet \wedge b_n = 0_n) \\(a_n \ll b_n)^\bullet &\triangleq (b_n^\bullet \wedge b_n \geq n)\end{aligned}$$

$$\begin{aligned}((a[i] \leftarrow e)[j])^\bullet &\triangleq (\text{ite } (i = j) \ e \ (a[j]))^\bullet \\&\triangleq ((i = j)^\bullet \wedge (\text{ite } (i = j) \ e^\bullet \ (a[j])^\bullet)) \\&\quad \vee (e^\bullet \wedge (a[j])^\bullet \wedge (e = a[j])) \\&\triangleq (i^\bullet \wedge j^\bullet \wedge (\text{ite } (i = j) \ e^\bullet \ (a[j])^\bullet)) \\&\quad \vee (e^\bullet \wedge (a[j])^\bullet \wedge (e = a[j]))\end{aligned}$$

Boolector: an efficient QF-solver for bitvectors and arrays

Best approaches

		Z3	Btor [•]	Btor [•] ▷ Z3
SMT-LIB	SAT	261	399	485
	# UNSAT	165	N/A	165
	UNKNOWN	843	870	619
	total time	270 150	350	94 610
BINSEC	SAT	953	1042	1067
	# UNSAT	319	N/A	319
	UNKNOWN	149	379	35
	total time	64 761	1 152	1 169

GRUB example

	Z3	Btor [•]
SAT	1	540
# UNSAT	42	N/A
UNKNOWN	852	355
total time	159 765	16 732

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]
	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

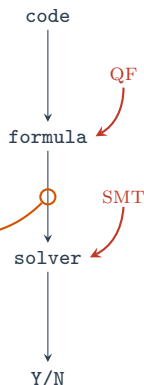
- Array theory useful for modelling memory or data structures...
- ...but a bottleneck for resolution of large formulas (BMC, SE)

- Existing approaches

- General decision procedures for the theory of arrays
- Dedicated handling of arrays inside tools

- Our proposal

- FAS, an efficient **simplification** for array theory
- ⇒ Improves existing solvers



Published in Logic for Programming, Artificial Intelligence and Reasoning, Awassa, Ethiopia, 2018 [LPAR18]

Presented in Journées Francophones des Langages Applicatifs, 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 $i \in \mathcal{I}$: $a[i] \leftarrow e$

$$\begin{aligned} \cdot [\cdot] &: \text{Array } \mathcal{I} \ \mathcal{E} \rightarrow \mathcal{I} \rightarrow \mathcal{E} \\ \cdot [\cdot] \leftarrow \cdot &: \text{Array } \mathcal{I} \ \mathcal{E} \rightarrow \mathcal{I} \rightarrow \mathcal{E} \rightarrow \text{Array } \mathcal{I} \ \mathcal{E} \end{aligned}$$

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

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

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```

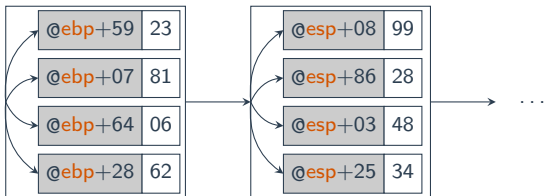


```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
assert (mem0 [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 $mem_1 [esp_1 + 48] = 1415$
 - \Rightarrow **scalability** issue: 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 $i=j$
 - 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...
- ⇒ Reduce the number of bases
-

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

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

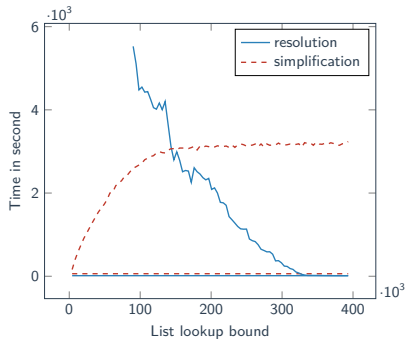
- 6,590 × 3 medium-size formulas from static SE
- TIMEOUT = 1,000 seconds

		simpl. time	#TIMEOUT and resolution time						#ROW
			Boolector		Yices		Z3		
concrete	default	61	0	163	2	69	0	872	866,155
	FAS	85	0	94	2	68	0	244	1,318
	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
symbolic	default	61	0	6,173	3	1,961	65	305,619	866,155
	FAS	91	0	6,117	3	1,965	66	158,635	531,654
	FAS-itv	111	0	4,767	2	1,108	43	80,569	295,333

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

		simpl. time	#TIMEOUT and resolution time						#ROW
			Boolector		Yices		Z3		
concrete	default	44	10	159	4	1,098	26	3.33	1,120,798
	FAS-list	1,108	8	845	4	198	10	918	456,915
	FAS	196	8	820	4	196	10	922	456,915
	FAS-itv	210	4	654	1	12	4	1,120	0
interval	default	44	12	131	12	596	27	0.19	1,120,798
	FAS-list	222	12	129	12	595	26	236	657,594
	FAS	231	12	129	12	597	26	291	657,594
	FAS-itv	237	12	58	12	28	19	81	651,449
symbolic	default	40	12	1,522	12	1,961	27	0.13	1,120,798
	FAS-list	187	11	1,199	12	2,018	26	486	657,594
	FAS	194	11	1,212	12	2,081	26	481	657,594
	FAS-itv	200	11	1,205	12	2,063	26	416	657,594

- Huge formula obtained from the ASPack packing tool
- 293 000 ROWS
- 24 hours of resolution!



Using FAS

- #ROW reduced to 2 467
- 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
 - 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
 - distinguish between controlled and uncontrolled inputs
 - 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
}

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;
}

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");
    }
}
```

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

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

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

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

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

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

i : input m : memory

s : secret p : stack pointer

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

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

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

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

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

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

Sending the formula to a solver:

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

- 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 \forall -formula is difficult
 - going low-level significantly increases formula size
- \Rightarrow The Death of SMT Solvers

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

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

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

$m_1 \triangleq m_0 [p_0 - \text{SIZE} \dots p_0 - 1] \leftarrow s$

$m_2 \triangleq m_1 [p_0 - 2 \cdot \text{SIZE} \dots p_0 - 2 \cdot \text{SIZE} + N - 1] \leftarrow i$

$i[0 \dots \text{SIZE} - 1] = i[\text{SIZE} \dots 2 \cdot \text{SIZE} - 1]$

$\wedge N \geq 2 \cdot \text{SIZE}$

Problems:

- finding a model for a \forall -formula is difficult [CAV18]
- going low-level significantly increases formula size [LPAR18]

\Rightarrow ~~The Death of SMT Solvers~~

For example with $\text{SIZE} = 8$,

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

- set of crackme challenges
- compare true and false positives

	SE robust		
	true positives	false positives	UNKNOWN
Boolector	N/A	N/A	24
CVC4	5	0	19
Yices	N/A	N/A	24
Z3	7	0	17

	SE classic		
	true positives	false positives	UNKNOWN
Boolector	12	11	1
CVC4	7	9	8
Yices	7	11	6
Z3	12	12	0

	SE robust + elim.		
	true positives	false positives	UNKNOWN
Boolector	12	0	12
CVC4	7	0	17
Yices	7	0	17
Z3	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:
solves `TIMEOUT`

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

Section 4

Conclusion

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

- ① Model Generation for Quantified Formulas
 - Proposed a novel and generic taint-based approach
 - Proved its correctness and its efficiency
 - Presented an implementation for arrays and bit-vectors
 - Evaluated on SMT-LIB and formulas generated by Symbolic Execution
- ② 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
- ③ Robust Symbolic Execution
 - Highlighted the problem of false positives in classic Symbolic Execution
 - Introduced formally the framework of Robust Symbolic Execution
 - Implemented a proof of concept in the binary analyser `BINSEC`

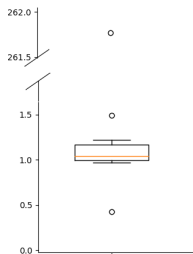


Guillaume Girol, B. Farinier, and S. Bardin.

Not all bugs are created equal, but robust reachability can tell the difference.
In *CAV 2021, Virtual, July 18-24, 2021*.

- Formal definition of Robust Reachability, application to SE and BMC
- Adaptation of standard optimizations to Robust Reachability
- Evaluation against 46 reachability problems including CVE replays and CTFs

	SE	BMC	RSE	RSE+	RBMC
Correct	30	22	37	44	32
False positive	16	14			
Inconclusive			7		1
Resource exhaustion		10	2	2	13



- Universal quantification of formulas has a cost, but not so high.
- RSE+ (robust SE with path merging) is 15% slower than SE in median, but with large outliers.