

# Black-box Code Analysis for Reverse Engineering

**Grégoire Menguy**, CEA LIST, France

Sébastien Bardin, CEA LIST, France

Nadjib Lazaar, LIRMM, France

Arnaud Gotlieb, Simula, Norway

Richard Bonichon, Tweag IO, France

Cauim de Souza Lima, CEA LIST, France



# Speaker



## Grégoire Menguy



PhD student at CEA LIST @BinsecTool



@grmenguy



<https://gregoiremenguy.github.io/>

# Context

Software size  +



 Hard to verify code



 Hard to test code



 Hard to understand code

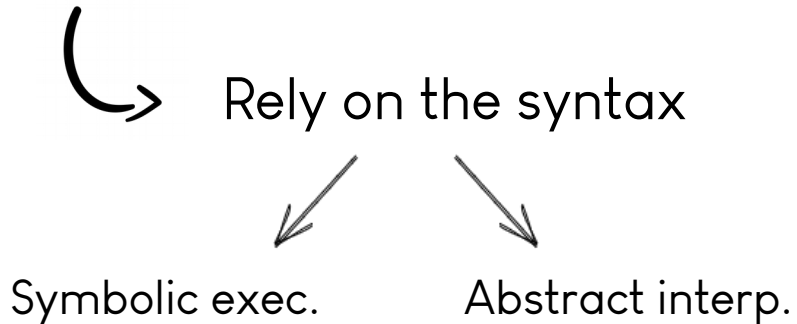


Need for automatic software analysis



# White vs Black-box

## White-box analysis



↳ Scale issues

## Black-box analysis

- ↳ Rely on code executions
- ↳ Insensitive to syntactic complexity



# Scenarios

## Deobfuscation

- ↳ Malware analysis
- ↳ Assess obfuscation strength



## Contract inference

- ↳ Core refactoring
- ↳ Code understanding

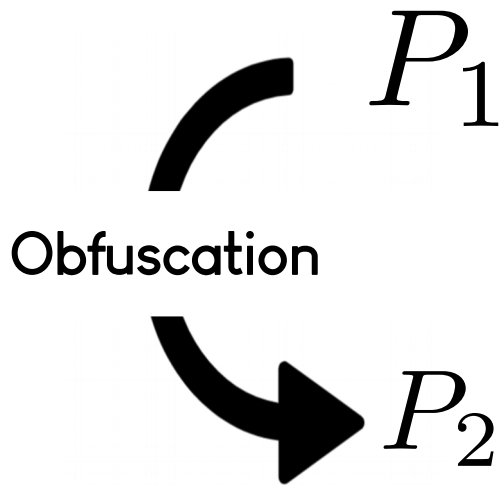


# Search-based Local Blackbox Deobfuscation: Understand, Improve and Mitigate

**Grégoire Menguy**, CEA LIST, France  
Sébastien Bardin, CEA LIST, France  
Richard Bonichon, Tweag IO, France  
Cauim de Souza Lima, CEA LIST, France

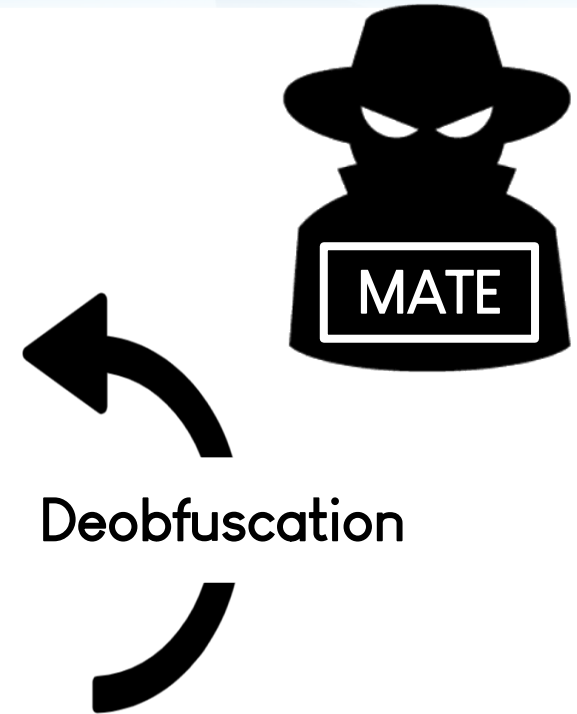


# Obfuscation



```
int f(in * l);  
int main();
```

```
double L,o,P,  
=dt,T,Z,D=1,d,  
s[999],E,h= 8,  
I,J,K,w[999],M,  
m,0,n[999],j=
```



# Deobfuscation

## Protecting Software through Obfuscation: Can It Keep Pace with Progress in Code Analysis?

SEBASTIAN SCHRITTWIESER, St. Pölten University of Applied Sciences, Austria  
STEFAN KATZENBEISSER, Technische Universität Darmstadt, Germany  
JOHANNES KINDER, Royal Holloway, University of London, United Kingdom  
GEORG MERZDOVNIK and EDGAR WEIPPL, SBA Research, Vienna, Austria

## A Generic Approach to Automatic Deobfuscation of Executable Code

Babak Yadegari    Brian Johannesmeyer    Benjamin Whitely    Saumya Debray  
Department of Computer Science  
The University of Arizona  
Tucson, AZ 85721  
{babaky, bjohannesmeyer, whitely, debray}@cs.arizona.edu

## Symbolic deobfuscation: from virtualized code back to the original\*

Jonathan Salwan<sup>1</sup>, Sébastien Bardin<sup>2</sup>, and Marie-Laure Potet<sup>3</sup>

## Backward-Bounded DSE: Targeting Infeasibility Questions on Obfuscated Codes\*

Sébastien Bardin  
CEA, LIST,  
91191 Gif-Sur-Yvette, France  
sebastien.bardint@cea.fr

Robin David  
CEA, LIST,  
91191 Gif-Sur-Yvette, France  
robin.david@cea.fr

Jean-Yves Marion  
Université de Lorraine,  
CNRS and Inria, LORIA, France  
jean-yves.marion@loria.fr



# BINSEC

# TRILON

Dynamic Binary Analysis





# Deobfuscation

## Protecting Software through Obfuscation: Can It Keep Pace with Progress in Code Analysis?

SEBASTIAN SCHRITTWIESER, St. Pölten University of Applied Sciences, Austria  
STEFAN KATZENBEISSER, Technische Universität Darmstadt  
JOHANNES KINDER, Royal Holloway, University of London  
GEORG MERZDOVNIK and EDGAR WEIPPL, SB

## Backward-Bounded DSE: Targeting Infeasibility Questions on Obfuscated Codes\*

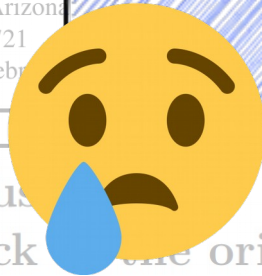
Robin David  
CEA, LIST,  
Furberville, France  
robin.david@cea.fr

Jean-Yves Marion  
Université de Lorraine,  
CNRS and Inria, LORIA, France  
jean-yves.marion@loria.fr

## A Generic Approach to Automatic Deobfuscation

Babak Yadegari    Brian Johannesmeyer    Benjamin  
Department of Computer Science  
The University of Arizona  
Tucson, AZ 85721  
{babaky, bjohannesmeyer, whitely, deb

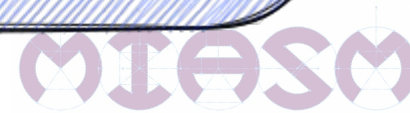
White-box deobfuscation  
is highly efficient



## Symbolic deobfuscation from virtualized code back to the original\*

Jonathan Salwan<sup>1</sup>, Sébastien Bardin<sup>2</sup>, and Marie-Laure Potet<sup>3</sup>

SEC  
TRILION  
Dynamic Binary Analysis



# Anti-White-Box Deobfuscation

But efficient countermeasures

## Information Hiding in Software with Mixed Boolean-Arithmetic Transforms

Yongxin Zhou, Alec Main, Yuan X. Gu, and Harold Johnson

Cloakware Inc., USA

{yongxin.zhou,alec.main,yuan.gu,harold.johnson}@cloakware.com



## How to Kill Symbolic Deobfuscation for Free (or: Unleashing the Potential of Path-Oriented Protections)

Mathilde Ollivier  
CEA, LIST,  
Paris-Saclay, France  
mathilde.ollivier2@cea.fr

Richard Bonichon  
CEA, LIST,  
Paris-Saclay, France  
richard.bonichon@cea.fr

Sébastien Bardin  
CEA, LIST,  
Paris-Saclay, France  
sebastien.bardin@cea.fr

Jean-Yves Marion  
Université de Lorraine, CNRS, LORIA  
Nancy, France  
Jean-Yves.Marion@loria.fr

## Probabilistic Obfuscation through Covert Channels

Jon Stephens   Babak Yadegari   Christian Collberg   Saumya Debray   Carlos Scheidegger


*Department of Computer Science*

*The University of Arizona*

*Tucson, AZ 85721, USA*

*Email: {stephensj2, babaky, collberg, debray, cscheid}@cs.arizona.edu*

# New Threat: Black-box Deobfuscation



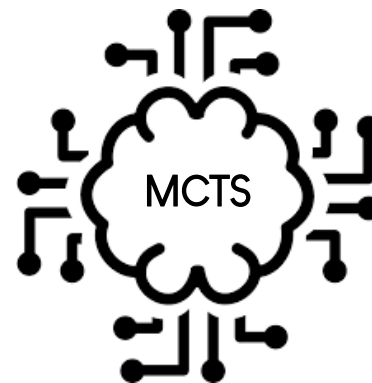
**Syntia: Synthesizing the Semantics of Obfuscated Code**

Tim Blazytko, Moritz Contag, Cornelius Aschermann,  
and Thorsten Holz, *Ruhr-Universität Bochum*

<https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/blazytko>

This paper is included in the Proceedings of the  
**26th USENIX Security Symposium**  
August 16–18, 2017 • Vancouver, BC, Canada  
ISBN 978-1-931971-40-9

Open access to the Proceedings of the  
26th USENIX Security Symposium  
is sponsored by USENIX



Bypasses white-box  
methods limitations

# Open Questions

## Understand



- Strengths ?
- Weaknesses ?
- Why ?

## Improve



- Why MCTS ?
- Can be improved ?
- Impacted by SOTA protections ?

## Mitigate



- How to protect ?

# Contributions

## Understand



- ➔ Propose missing formalisation
- ➔ Refine Syntia Xps: new strengths & weaknesses
- ➔ Show & explain why MCTS not appropriate

## Improve



- ➔ S-metaheuristics > MCTS
- ➔ Implement our approach: Xyntia
- ➔ Evaluation of Xyntia

## Mitigate



- ➔ Propose 2 protections
- ➔ Evaluate them against Xyntia and Syntia



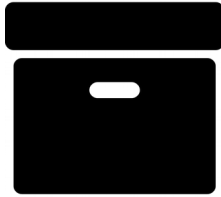
# Black-box deobfuscation

## What's that ?

# Black-box Deobfuscation

## 1. Sample

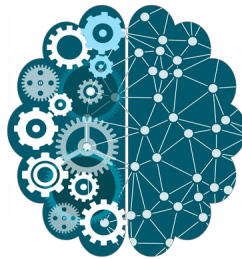
$(t = 1, T = 2)$   
 $(t = 2, T = 5)$   
 $(t = 0, T = 6)$   
...



-1  
-3  
-6  
...

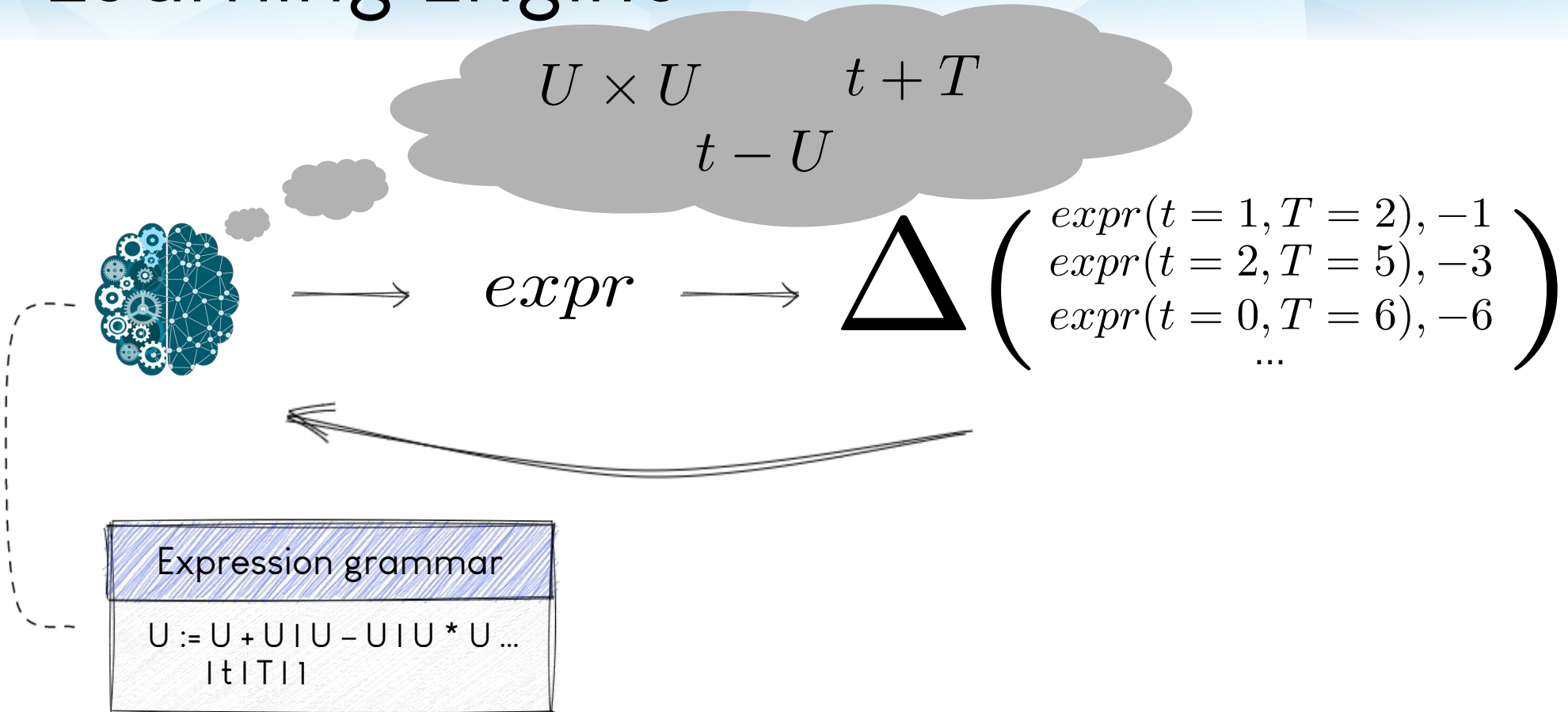
## 2. Learn

$(t = 1, T = 2) \rightarrow -1$   
 $(t = 2, T = 5) \rightarrow -3$   
 $(t = 0, T = 6) \rightarrow -6$   
...



$t - T$

# Learning Engine





# Why Black-box ?

Given a language  $L$  and an expression “ $e$ ” in  $L$

Syntactic complexity

Size of the expression “ $e$ ”

Semantic complexity

Size of the smallest expr. in  $L$   
equivalent to “ $e$ ”

Example

$t - T$  is syntactically simpler than  $(t \vee -2T) \times 2 - (t \oplus -2T) + T$

**but** they share the same semantic complexity (being equivalent)

# Why Black-box ?

Given a language  $L$  and an expression “ $e$ ” in  $L$

Syntactic complexity

Size of the expression “ $e$ ”

Semantic complexity

Size of the smallest expr. in  $L$   
equivalent to “ $e$ ”

Example

$t - T$  is syntactically simpler than  $(t \vee -2T) \times 2 - (t \oplus -2T) + T$

but they share the same semantic complexity (being equivalent)

Obfuscation increases syntactic complexity  
→ No impact on black-box methods

Understand



# Zoom of SOTA: Syntia


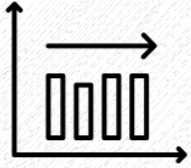


**Dig into Syntia and deepen its evaluation**


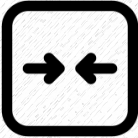
- RQ1: stability of Syntia
- **RQ2: efficiency of Syntia**
- RQ3: impact of operators set

# Syntia: New Results


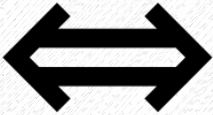
Stable



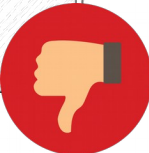

Quality



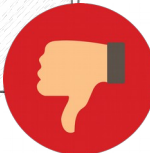

Correctness



Speed



Robustness



# Experimental Design

B1 (Syntia)
<ul style="list-style-type: none"><li>- 500 expressions</li><li>- Up to 3 inputs</li><li>- redundancy</li><li>- Unbalances w.r.t. type</li></ul>

B2 (ours)
<ul style="list-style-type: none"><li>- 1110 expressions</li><li>- 2 to 6 inputs</li><li>- No redundancy</li><li>- Balances w.r.t. type</li></ul>

	Type			# Inputs				
	Bool.	Arith.	MBA	2	3	4	5	6
#Expr.	370	370	370	150	600	180	90	90

**Table 1: Distribution of samples in benchmark B2**

# Evaluation of Syntia

## B1 (Syntia)

- With a 60 s/expr. timeout: 75% of success rate
- With a 1 h/expr. timeout: 88% of success rate
- With a 12 h/expr. timeout: **97% of success rate**

## B2 (Ours)

**Table 2: Syntia depending on the timeout per expression (B2)**

	1s	10s	60s	600s
Succ. Rate	16.5%	25.6%	34.5%	42.3%
Equiv. Range	16.3%	25.1 - 25.3%	33.7 - 34.0%	41.4 - 41.6%
Mean Qual	0.35	0.49	0.59	0.67

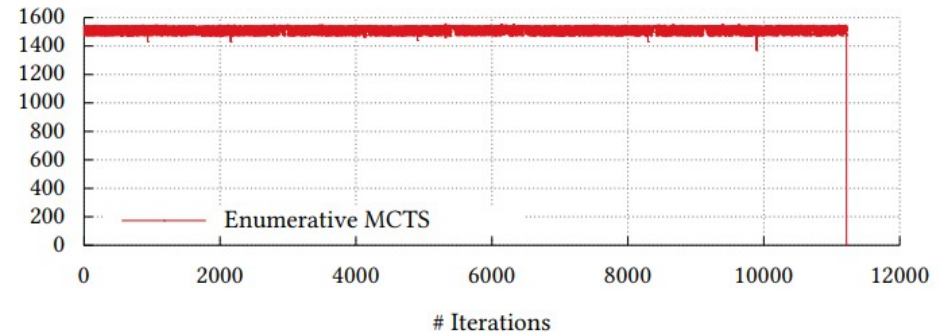
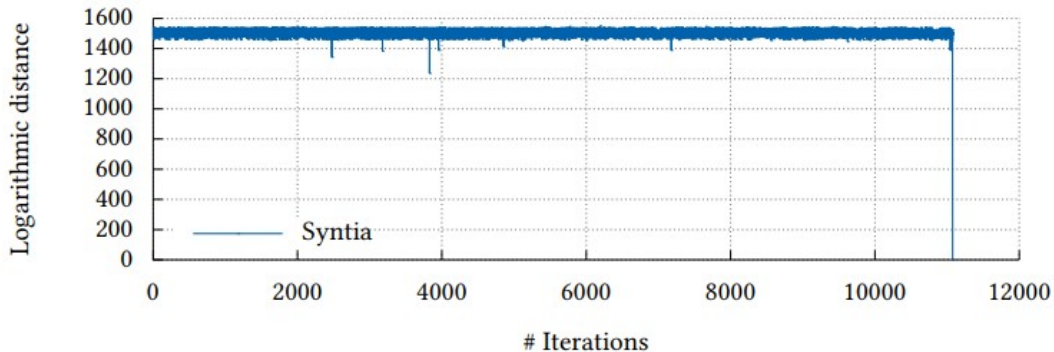
# Why ?

- Syntia manipulates non terminal expressions, e.g.,  $U - V$
- Scoring of non terminal expressions can be misleading



$$U - V \rightsquigarrow \begin{cases} t - T \\ t - 1 \\ 1 - 1 \end{cases}$$

- Syntia (i.e., MCTS) = “almost BFS”



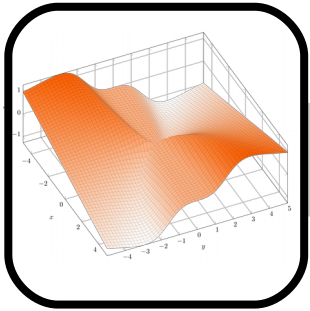


Improve



# Black-box Deobf: An Optimization Pb

Syntia sees blackbox deobfuscation as a **single player game**

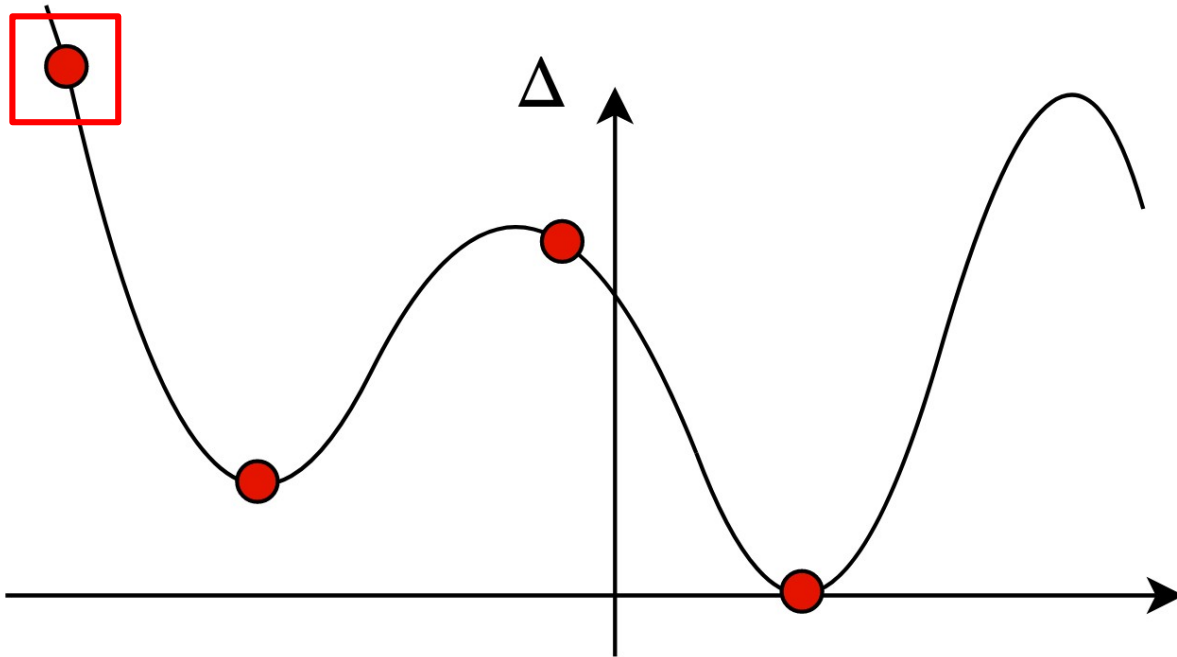


We propose to see it as an **optimization problem**

↪ **Goal** : find  $\underbrace{s^*}_{\text{An expr.}}$  s.t.  $\underbrace{f(s^*)}_{\Delta} \leq f(s), \forall s \in S$

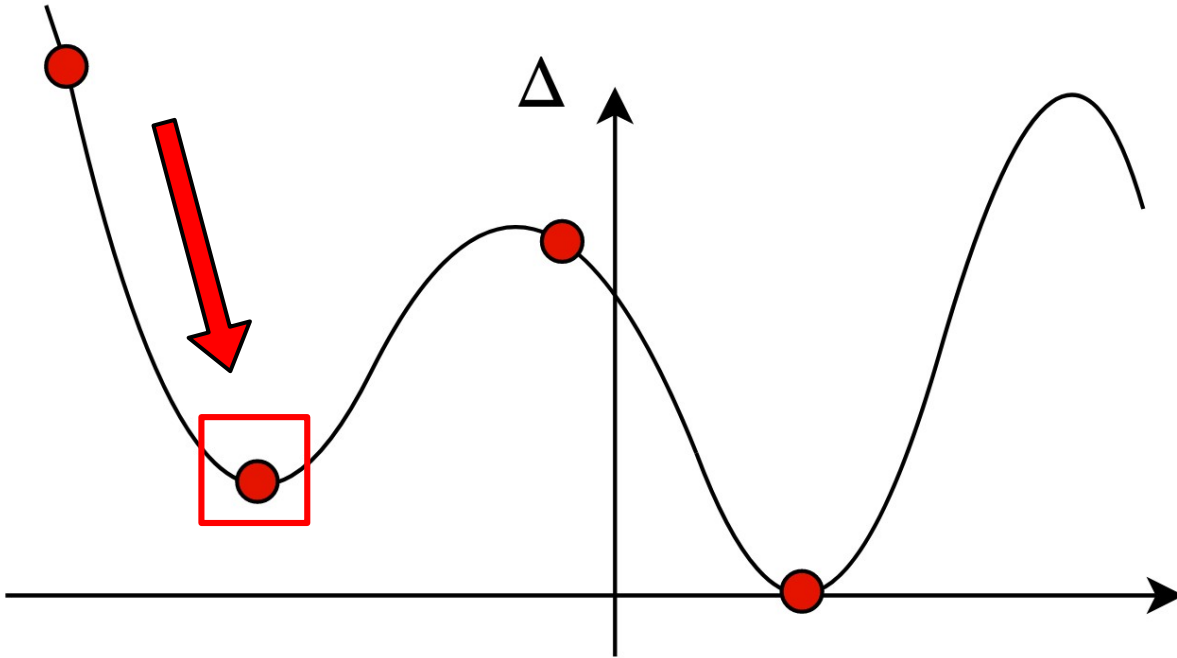
# S-metaheuristics

- Solve optimization problems



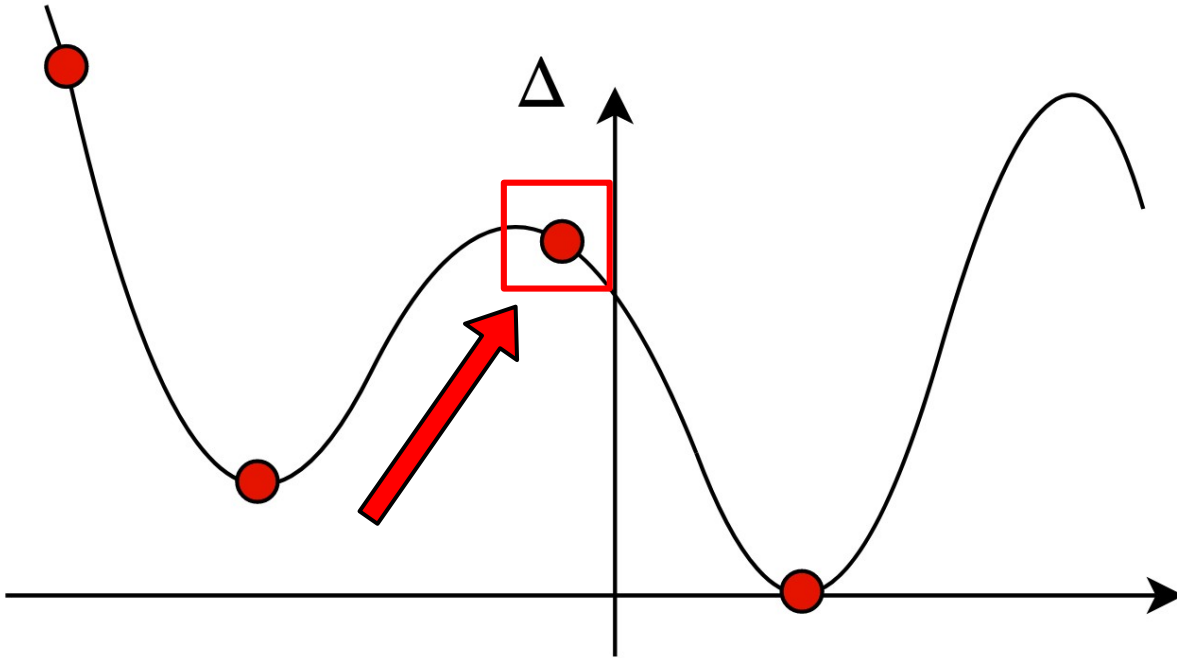
# S-metaheuristics

- Solve optimization problems



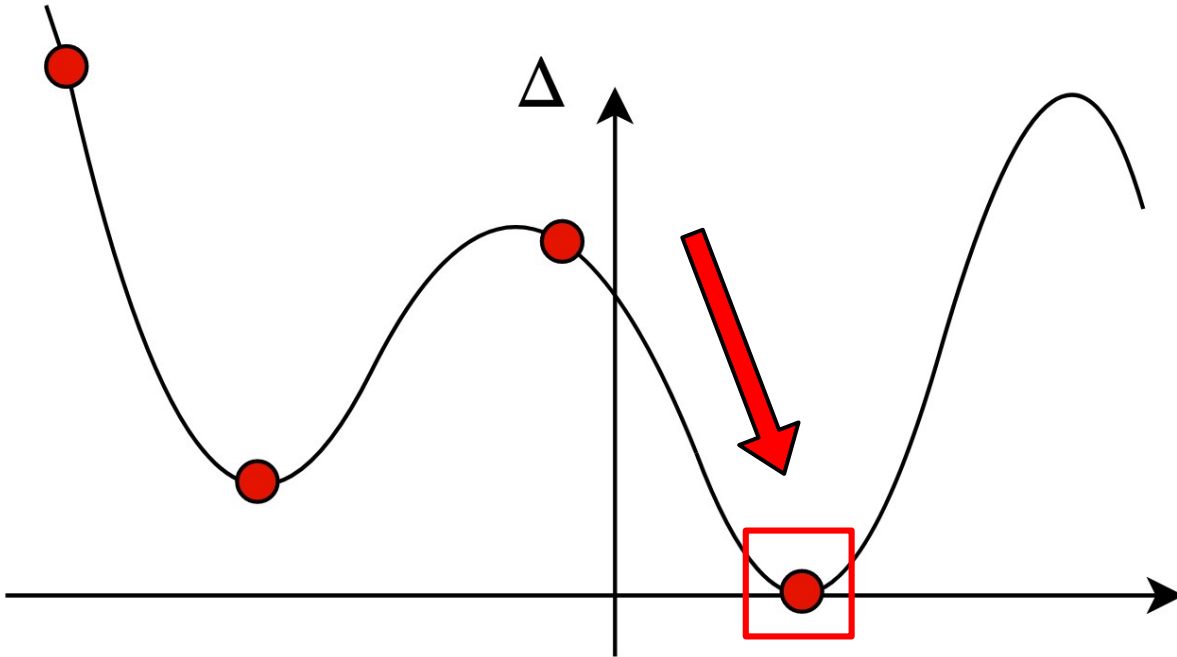
# S-metaheuristics

- Solve optimization problems



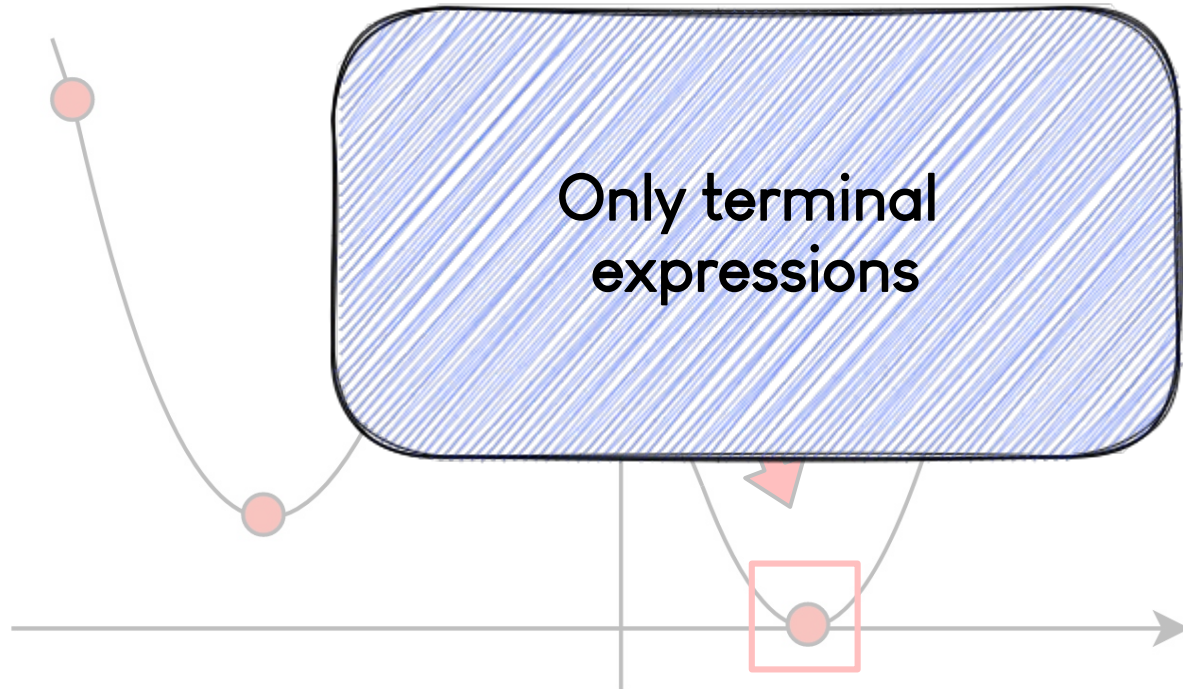
# S-metaheuristics

- Solve optimization problems

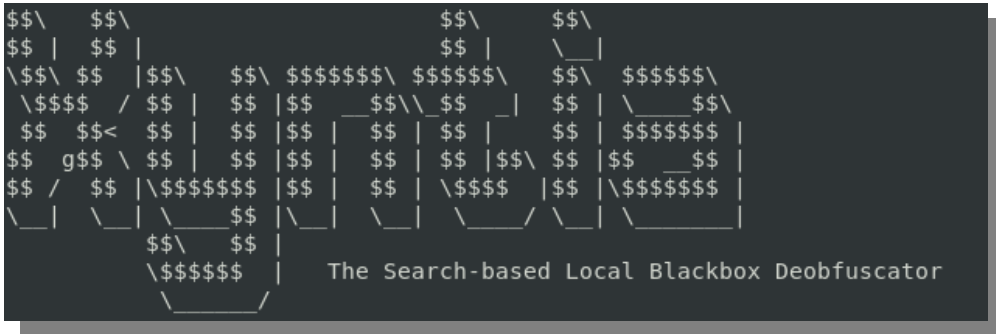


# S-metaheuristics

- Solve optimization problems

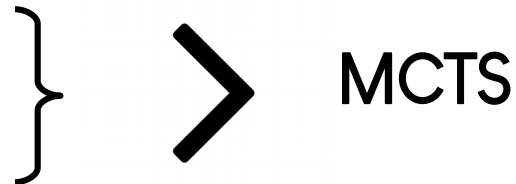


# New Prototype: Xyntia



↳ S-metaheuristics

- ↳ Can choose between:
- Hill Climbing
  - Simulated Annealing
  - Metropolis Hasting
  - Iterated Local Search





# Xyntia vs Syntia

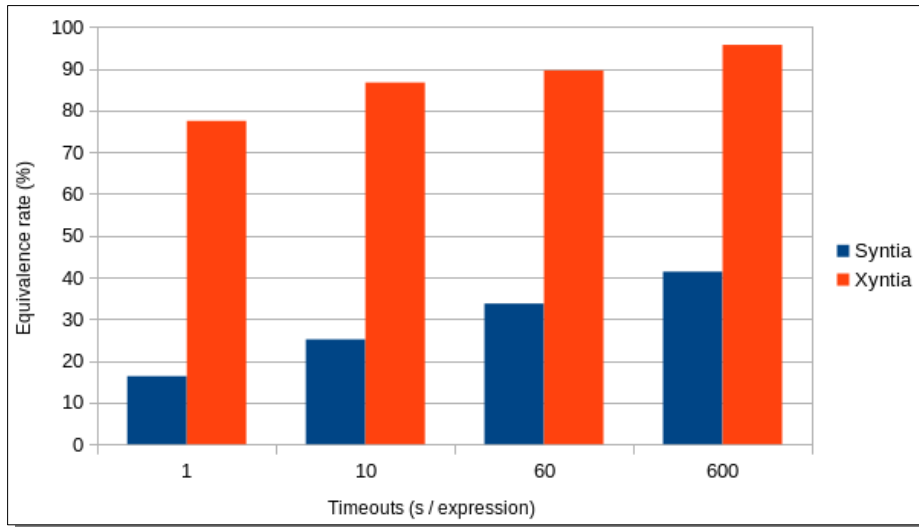
## B1 (Syntia)

– 100 % success rate in 1s/expr.



Syntia: 75% in 60 s/expr.

## B2 (Ours)



# Xyntia vs Syntia

BI (Syntia)

- 100 % success rate in 1s/expr



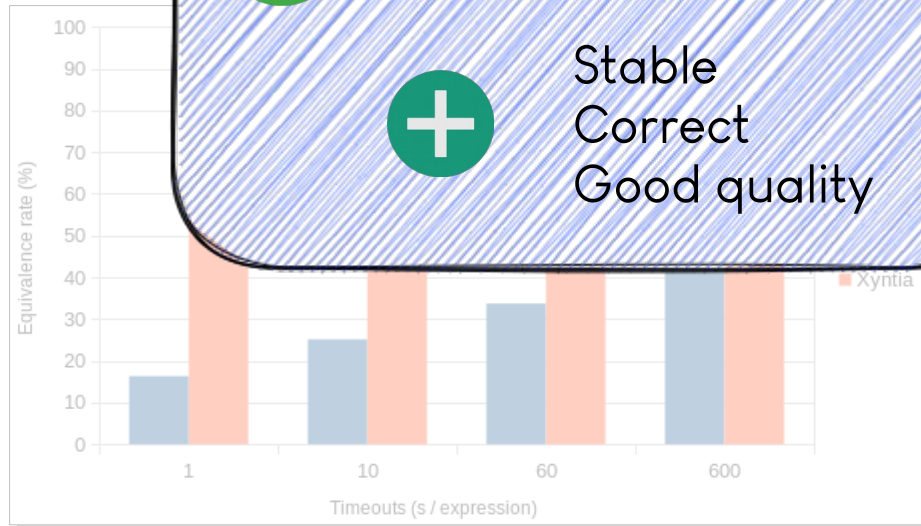
Syntia: 75% in 60 s/expr.



**Robust & Fast**

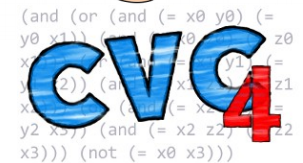


Stable  
Correct  
Good quality



# Other Experiments

- Xyntia vs Qsynth
- Xyntia vs “compiler like simplifications”
- Xyntia vs program synthesizer CVC4
- Xyntia vs superoptimizer STOKE
- Use-cases
  - ↳ State-of-the-art protections
  - ↳ VM-based obfuscation



404

Not Found

The resource requested could not be found on this server!



# What's Next ?



Mitigate



# Context: Virtualization

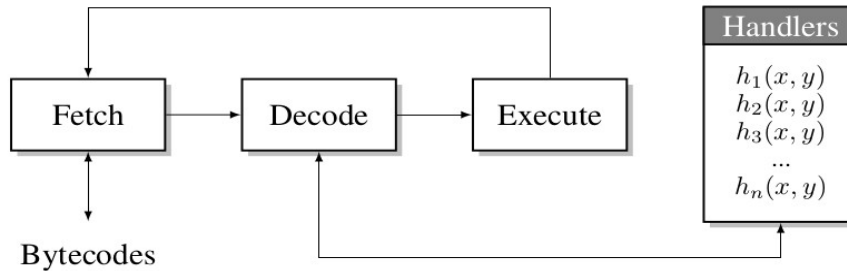


Fig. 1: Virtualization based obfuscation

Proved to be sensitive to black-box deobfuscation



**Themida**<sup>®</sup>  
ADVANCED WINDOWS SOFTWARE PROTECTION



# Why VM-based Obf. Is Vulnerable ?



- Handlers are too semantically simple
  - ↳ e.g.,  $+$ ,  $-$ ,  $\times$ ,  $\wedge$ ,  $\vee$
- Obfuscation increases syntactic complexity
  - ↳ Black-box deobf. is not impacted

We need to move ...

From syntactic to semantic complexity

# Semantically Complex Expressions

## Goal

- ↳ Increase the semantic complexity of each handlers
- ↳ Keep a Turing complete set of handlers

## Example

$$\begin{array}{l} h_0 = (x + y) + -((a - x^2) - (xy)) \\ + h_1 = (a - x^2) - xy + -(y - (a \wedge x)) \times (y \otimes x) \\ + h_2 = (y - (a \wedge x)) \times (y \otimes x) \\ \hline h = x + y \end{array}$$



# Merged Handlers

**Goal :** Increase handlers semantic complexity + sampling harder

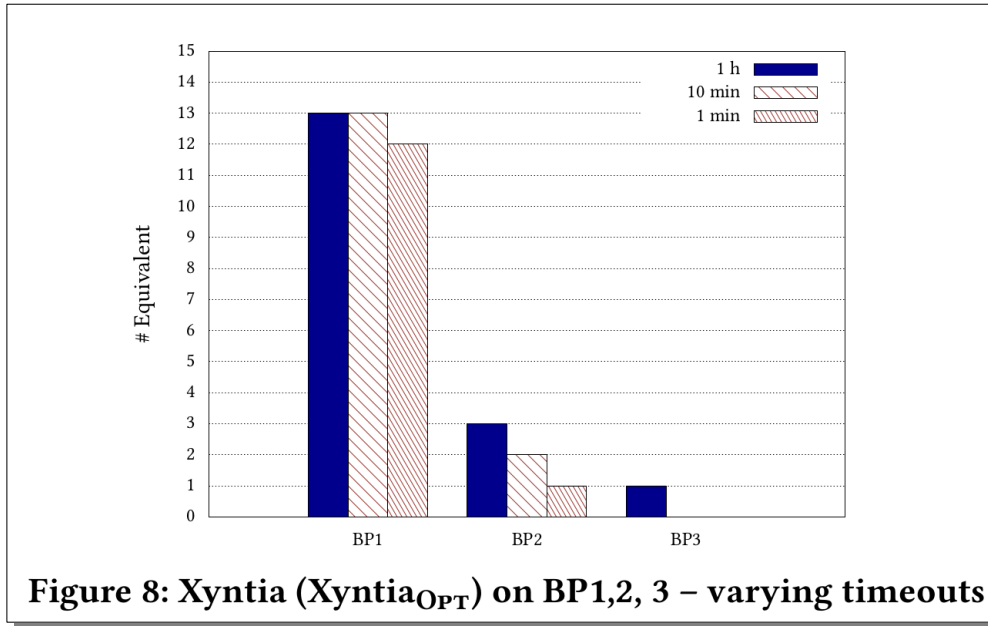
**Example :**  $h_1(x, y) = x + y$  and  $h_2(x, y) = x \wedge y$

—►  $h(x, y, c) = \text{if } (c = cst) \text{ then } h_1(x, y) \text{ else } h_2(x, y)$

**Hide conditionals :**

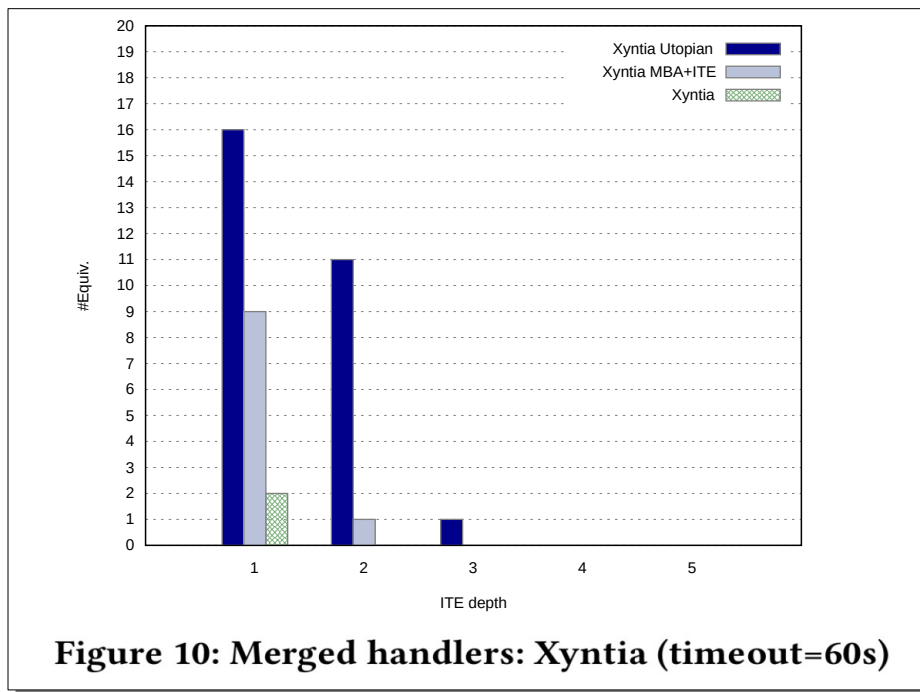
```
int32_t h(int32_t a, int32_t b, int32_t c) {  
    // if (c == cst) then h1(a,b,c) else h2(a,b,c);  
    int32_t res = c - cst ;  
    int32_t s = res >> 31;  
    res = (-((res ^ s) -s) >> 31) & 1;  
    return h1(a, b, c)*(1 - res) + res*h2(a, b, c);  
}
```

# Semantically Complex Handlers: Results



**More results:** Syntia with 12h/exprs. → 1/15 on BP1

# Merged Handlers: Results



**More results:** Syntia finds nothing for  $\geq 2$  nested ITE

# Xyntia: The Recap



**MCTS is not appropriate for blackbox deobfuscation**

→ Search space too unstable

→ Estimation of non terminal expressions pertinence is misleading



**S-metaheuristics yields a significant improvement**

→ More robust

→ Much Faster



**Moving for syntactic to semantic complexity**

→ 2 efficient methods to protect against blackbox deobfuscation

# Discussion: No Guarantees

– Xyntia & Syntia have **no correctness guarantees**

– But **some contexts need it:**

↳ Code verification

↳ Code refactoring



Can black-box approaches have clear guarantees ?



# Automated Program Analysis: Revisiting Precondition Inference through Constraint Acquisition

**Grégoire Menguy**, CEA LIST, France

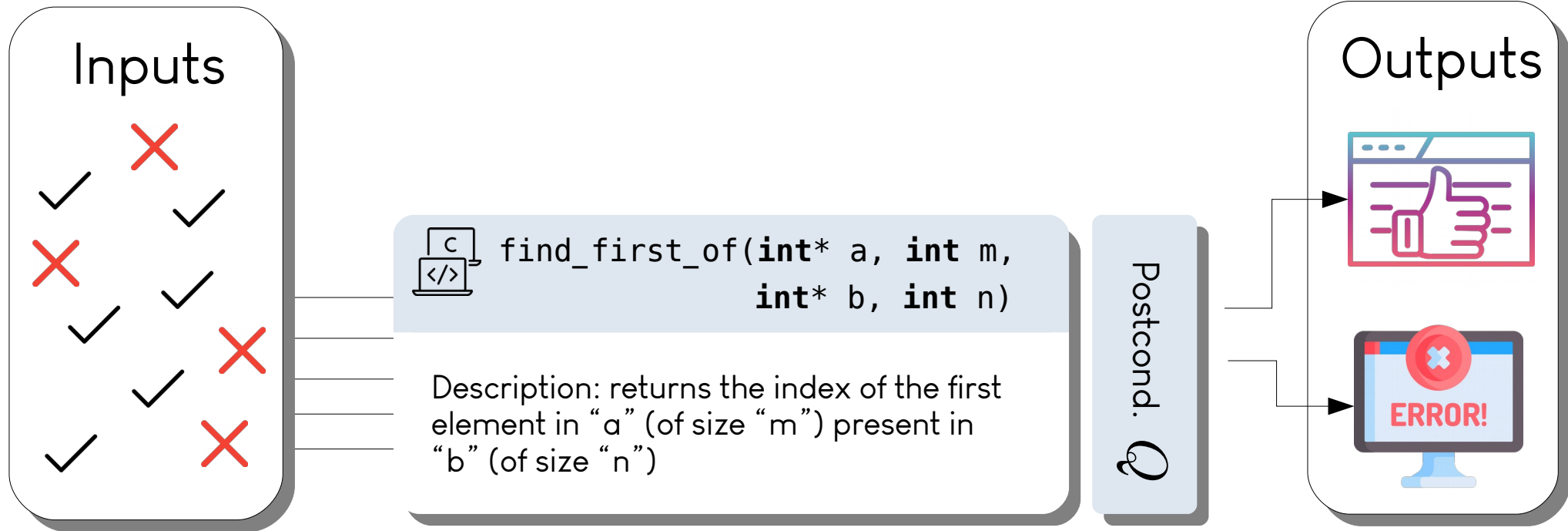
Sébastien Bardin, CEA LIST, France

Nadjib Lazaar, LIRMM, France

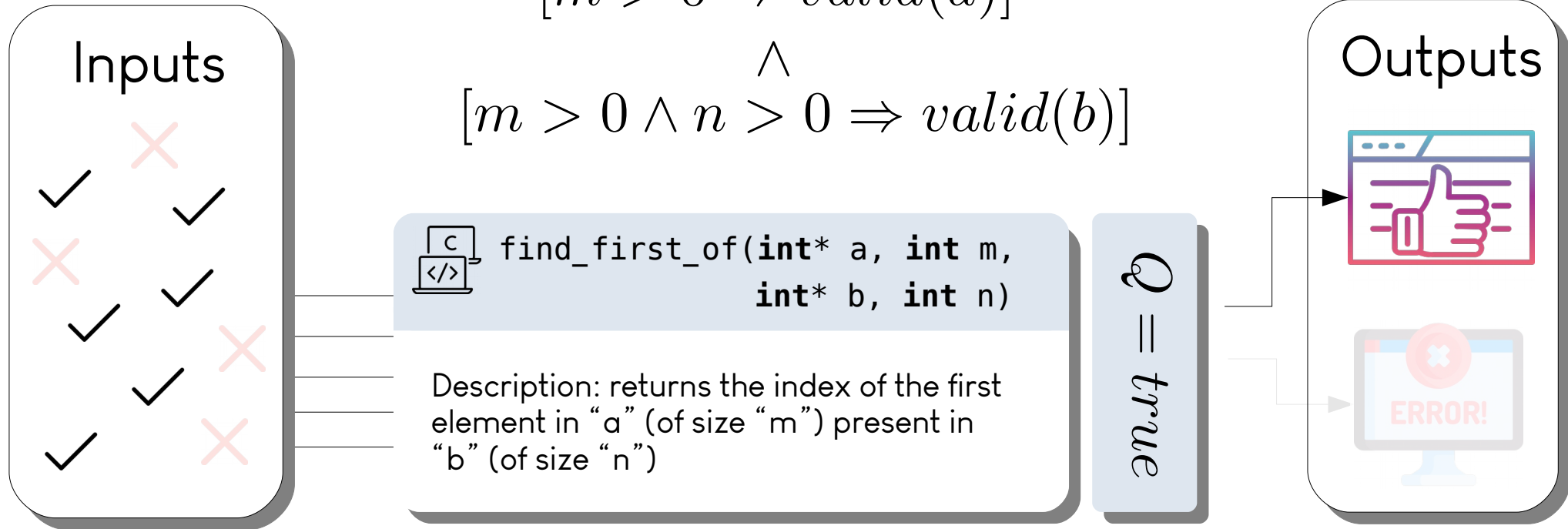
Arnaud Gotlieb, Simula, Norway



# Dream: Infer Preconditions



# Dream: Infer Preconditions

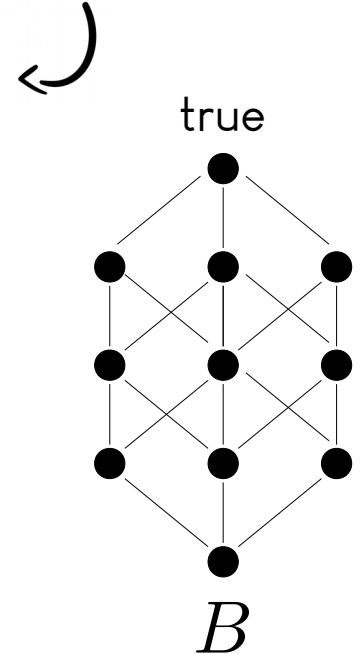
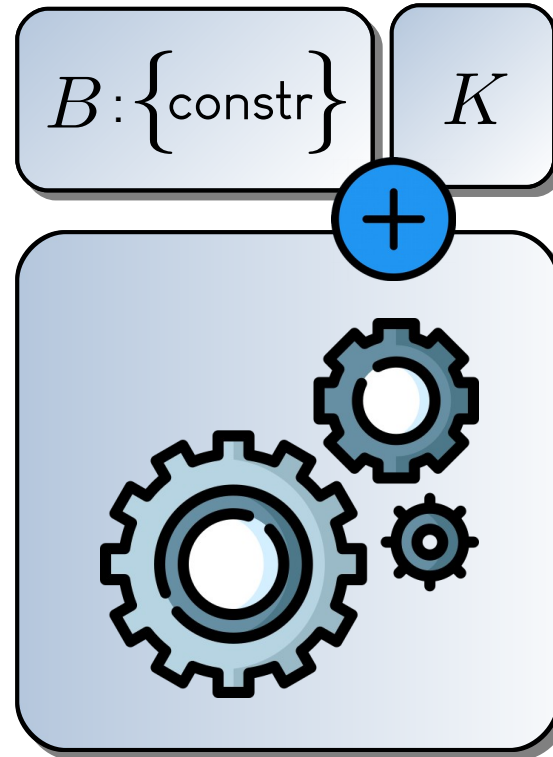


Undecidable problem: Rice theorem (1953)



# Active Constraint Acquisition

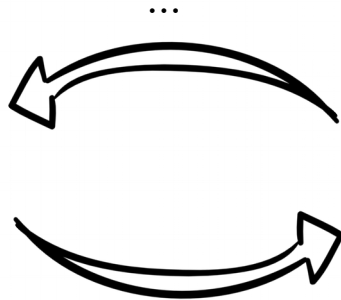
Background knowledge:  
rules to speed up learning



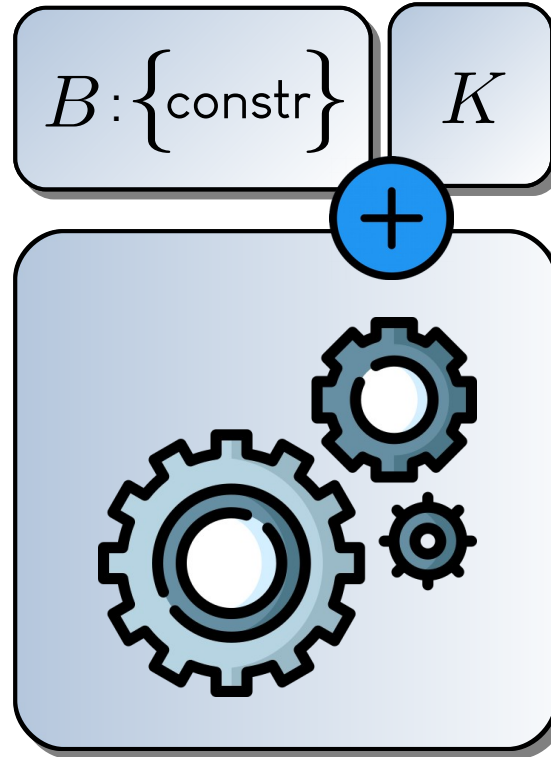
# Active Constraint Acquisition



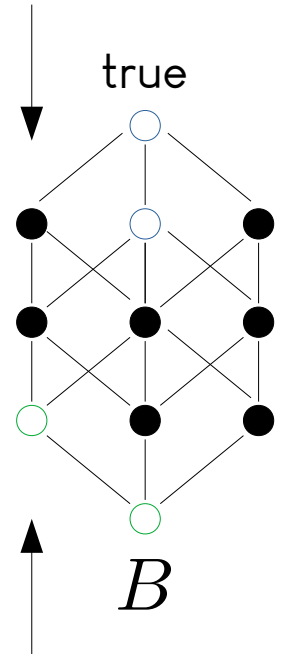
Query  
Elise: 8h - 12h  
Paul: 10h - 11h



yes /no

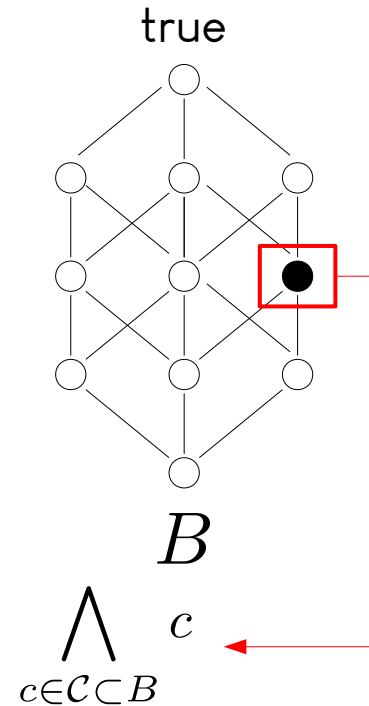
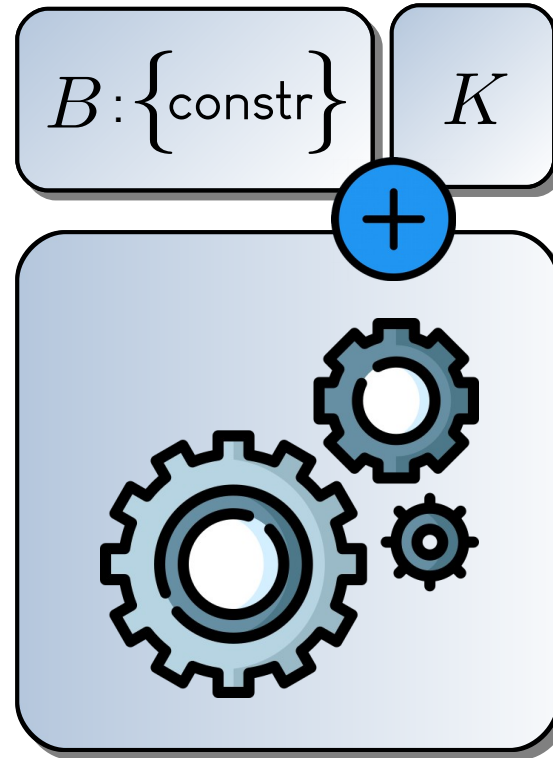
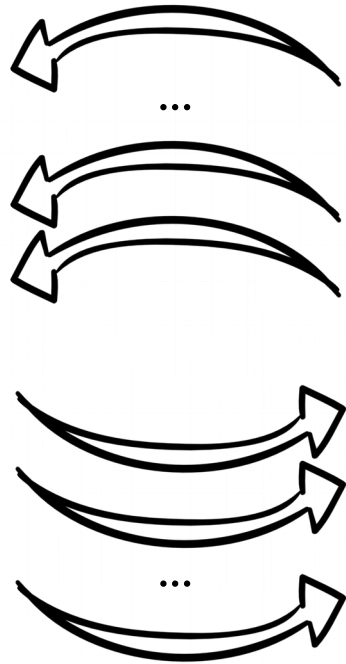


no: Top-down

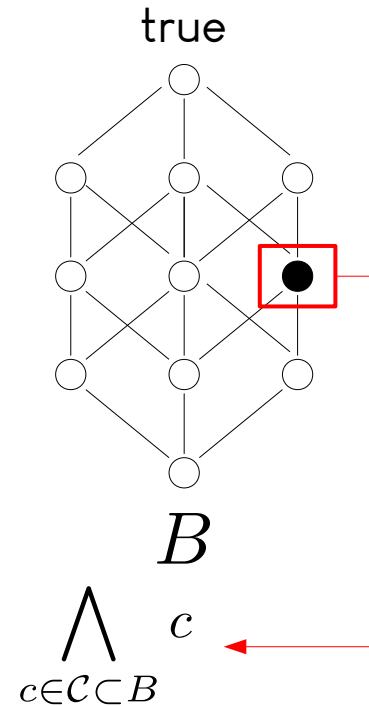
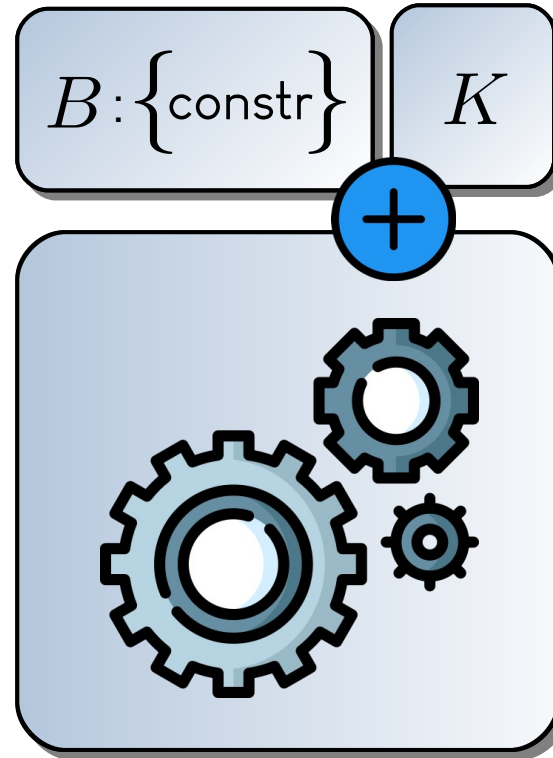


yes: Bottom-up

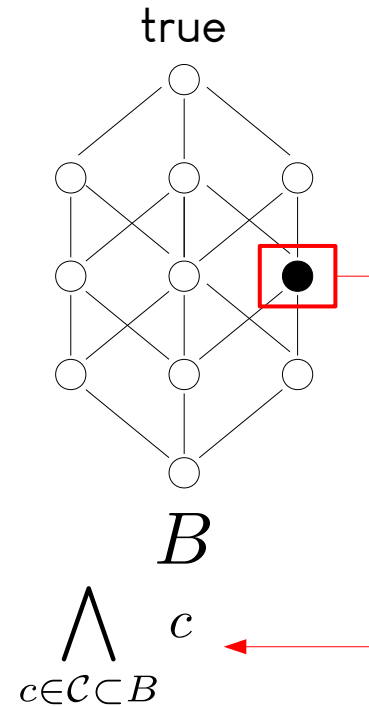
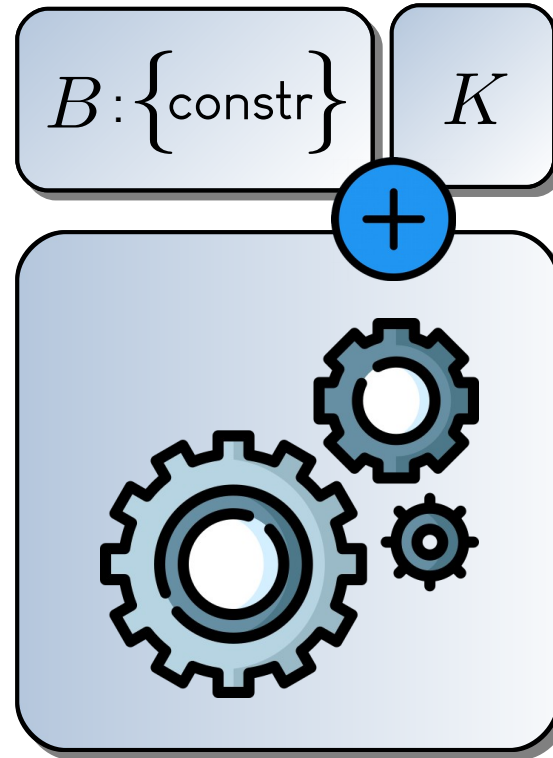
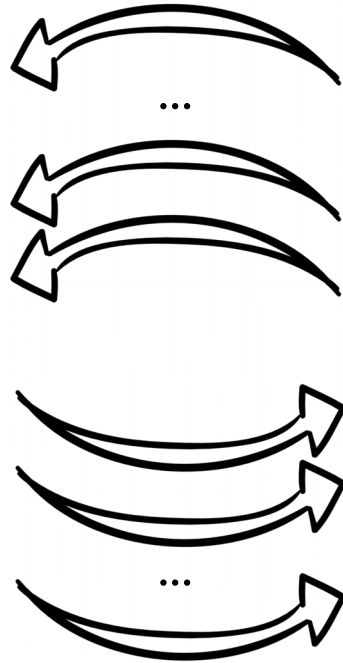
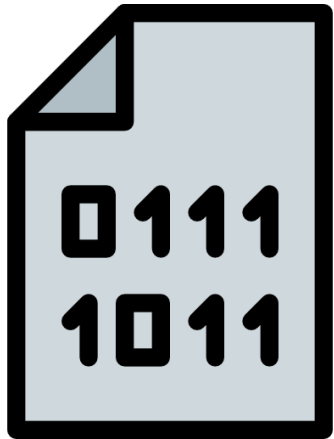
# Active Constraint Acquisition



# Careful: too many queries





# PreCA: Precondition Constr. Acq.



# Theoretical Analysis

PreCA guarantees

- ↳ If B is expressive enough .....  or Precond.
- ↳  If code exec. always terminates ..... The most general precondition

These are good theoretical guarantees

- ↳ SOTA executions based methods, from programming language community, have no clear guarantees

# Evaluation

**Dataset:** 94 learning tasks • compiled C functions (string.h, arrays, arithmetic ...)

**Evaluation:** \_\_\_\_\_

1 hour

PreCA

92%

41%



Daikon, PIE, Gehr et al

At most 52%

At most 23%



P-Gen

74%

34%



PreCA better in 5s than concurrent tools in 1 hour

# Conclusion

Black-box methods can be used in broad contexts

↳ Deobfuscation – Contracts inference

Balance between robustness and guarantees

← robustness correctness →





Thank you for your  
attention



@grmenguy



<https://gregoiremenguy.github.io/>