

DE LA RECHERCHE À L'INDUSTRIE



www.cea.fr

Sibyl

Function divination tool

Camille Mougey

- 1 Principe
 - Objectif
 - État de l'art
 - Idée
 - Implémentation

- 2 Démonstration

- 3 Ajout d'un test
 - Exemple : strlen

- 4 Focus sur l'apprentissage
 - Constat
 - Notre solution
 - Démonstration

- 5 Limites et futur

« Dans ce binaire / firmware / malware / shellcode / ..., la fonction à l'adresse 0x1234 est un memcpy »

Approche statique

- FLIRT
- Polichombr, Gorille, BASS
- Machine learning (ASM considéré comme NLP)
- *Bit-precise Symbolic Loop Mapping*
- ...

Approche dynamique / trace

- Entropie des données dans les boucles
- Avalanche de propagation de teinte

Approche proche de Sibyl

- Angr "identifier"^a \approx PoC pour le CGC

a. <https://github.com/angr/identifier>

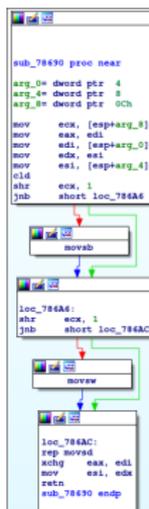


Figure – memcpy « naïf »

Problème

Comment reconnaître quand optimisé / vectorisé / compilateur différent / **obscurci** ?

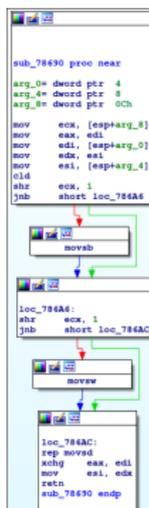


Figure – memcpy « naïf »

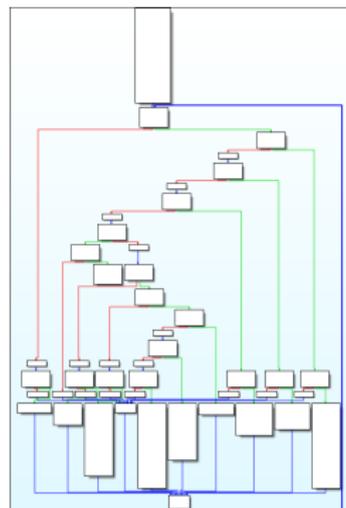


Figure – memcpy « obscurci »

Problème

Comment reconnaître quand optimisé / **vectorisé** / compilateur différent / obscurci ?

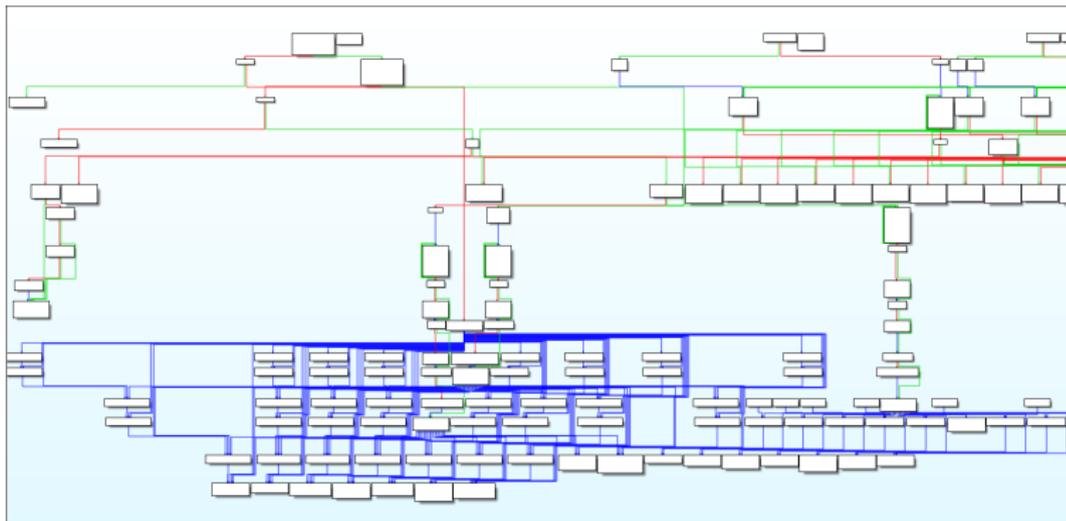


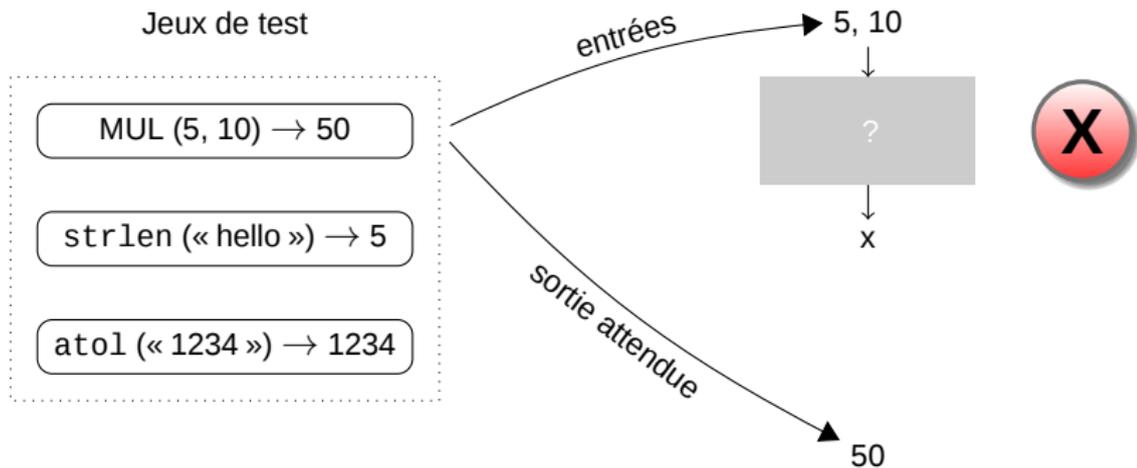
Figure – memcpy « SSE »

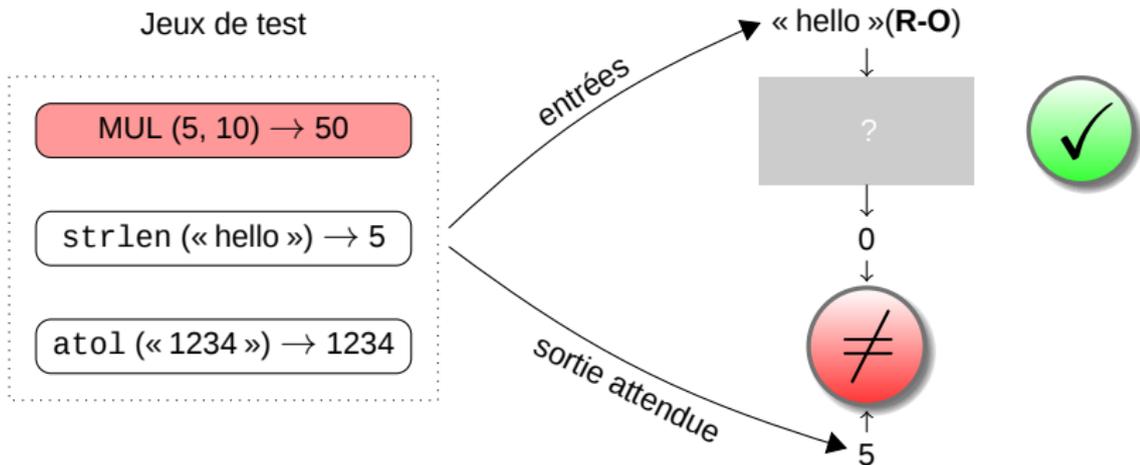
Idée

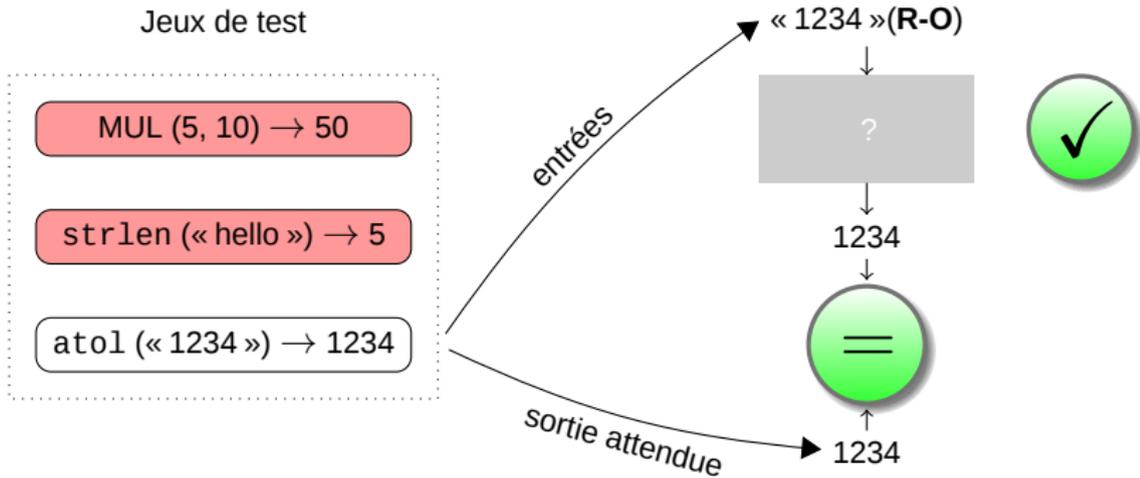
- Fonction = boîte noire
- Entrées choisies
- Sorties observées \leftrightarrow Sorties attendues

Spécifiquement

- Entrées = { arguments, mémoire initiale }
- Sorties = { valeur de retour, mémoire finale }
- Environnement minimaliste : { binaire mappé, stack }







Jeux de test

MUL (5, 10) → 50

strlen (« hello ») → 5

atol (« 1234 ») → 1234

atol

Le diable est dans les détails

- Fonctions attendant un *pointeur*, on donne l'entier 2
 - → Résistance aux crashes

Le diable est dans les détails

- Fonctions attendant un *pointeur*, on donne l'entier 2
 - → Résistance aux crashes
- Fonctions attendant un *compteur*, on donne un *pointeur*
 - → Détection des trèèèè longues boucles / temps d'exécution

Le diable est dans les détails

- Fonctions attendant un *pointeur*, on donne l'entier 2
 - → Résistance aux crashes
- Fonctions attendant un *compteur*, on donne un *pointeur*
 - → Détection des très très longues boucles / temps d'exécution
- Un appel peut ne pas être suffisant
 - (2, 2) → Func → 4
 - add, mul, pow?
 - → Politique de test : « test1 & (test2 || test3) »

Et tout le reste

- Pas d'interférence entre les tests
- *Embarrassingly parallel*
- Éviter les pertes de temps (instanciation de classe, stack d'appel, ...)
- Description de test « architecture agnostique », « ABI agnostique »
- ...

→ et si on en faisait un outil ?

Sibyl

- Open-source, GPL
- Version actuelle : 0.2
- CLI + Plugin IDA
- /doc
- D'autres contributeurs, notamment F. Desclaux & P. Graux



<https://github.com/cea-sec/Sibyl>

« Une sibylle est une « prophétesse », une femme qui fait œuvre de divination. »

Dependencies

```
$> sudo apt install python-pyparsing python-dev
```

Elfesteem

```
$> git clone github.com/serpilliere/elfesteem && cd elfesteem && \  
sudo python setup.py install
```

Miasm

```
$> git clone github.com/cea-sec/miasm && cd miasm && \  
sudo python setup.py install
```

Sibyl

```
$> git clone github.com/cea-sec/Sibyl && cd Sibyl && \  
sudo python setup.py install
```

Démonstration

- Présentation de la CLI
- Configuration
- Exemple : firmware
- Exemple : fonction inconnue optimisée, via IDA

Base de signature

- Située dans `sibyl/test/*`
- Voir `doc/ADD_TESTS.md`

strlen
strnicmp
strcpy
strncpy
strcat
strncat
strcmp
strchr
strrchr
strnlen
strspn
strpbrk
strtok

strsep
memset
memmove
stricmp
strrev
memcmp
bzero
strncmp
memcpy
abs
a64l
atoi
isalnum

isalpha
isascii
isblank
iscntrl
isdigit
isgraph
islower
isprint
ispunct
isspace
isupper
isxdigit

Création de la classe Python représentant le test

```
1 from sibyl.test.test import Test, TestSet
2
3 class TestStrlen(Test):
4     func = "strlen"
```

Préparation du premier test :

Allocation d'une chaîne de caractère en **read-only**

```

5      # Test1
6      my_string = "Hello, w%srld !"
7
8      def init(self):
9          # Alloue une chaine en READ-only
10         self.my_addr = self._alloc_string(self.my_string)
11
12         # Passe le pointeur dans le premier argument
13         self._add_arg(0, self.my_addr)
  
```

Vérification du résultat

```

14     def check(self) :
15         # Vérifie le résultat attendu
16         return self._get_result() == len(self.my_string)

```

Second test pour éviter les fonctions retournant toujours 15 (faux positif trivial)

```

17     # Test2 : évite les fonctions qui retourne toujours 15
18     def init2(self):
19         self.my_addr = self._alloc_string(self.my_string * 4)
20         self._add_arg(0, self.my_addr)
21
22     def check2(self):
23         return result == len(self.my_string * 4)

```

« Politique »de test

```
24      # Politique de test : test1 & test2  
25      tests = TestSetTest(init , check) & TestSetTest(init2 , check2)
```

```
1 class TestStrlen(Test):
2     func = "strlen"
3
4     # Test1
5     my_string = "Hello, w%srld !"
6
7     def init(self):
8         # Alloue une chaine en READ-only
9         self.my_addr = self._alloc_string(self.my_string)
10
11        # Passe le pointeur dans le premier argument
12        self._add_arg(0, self.my_addr)
13
14    def check(self):
15        # Vérifie le résultat attendu
16        return self._get_result() == len(self.my_string)
17
18    # Test2 : évite les fonctions qui retourne toujours 15
19    def init2(self):
20        self.my_addr = self._alloc_string(self.my_string * 4)
21        self._add_arg(0, self.my_addr)
22
23    def check2(self):
24        return result == len(self.my_string * 4)
25
26    # Politique de test : test1 & test2
27    tests = TestSetTest(init, check) & TestSetTest(init2, check2)
```

Ajout de nouveaux tests

- « à la main » pour chaque test
- assez rapide, mais ingrat
- compliqué à l'échelle d'une bibliothèque

→ Apprentissage automatique !

Version naive

- 1 Prendre un *snapshot* du contexte (mémoire + registres) *avant*
- 2 Faire tourner la fonction
- 3 Prendre un *snapshot* mémoire *après*

Mais ...

- « bruit » lié à l'implémentation
 - variable locale
 - constantes du binaire (par ex : SBOX AES)
- différent *packing* de structure
- changement d'architecture
- changement de convention d'appel (registre → stack)
- ...

Pré-requis

- Un binaire exécutant *au moins une fois* la fonction
 - Tests de regression
 - Malware
 - Binaire « maison » pour l'occasion

Pré-requis

- Un binaire exécutant *au moins une fois* la fonction
 - Tests de regression
 - Malware
 - Binaire « maison » pour l'occasion
- Signature *complète* de la fonction
 - Comprend le type des arguments, et structures associées
 - Obtenu des sources
 - Ou d'un reverse partiel de la fonction

Méthode

- 1 Trace d'exécution de la fonction
- 2 Élagage des runs trop semblables (*branch coverage*)
- 3 Rejeu en *Dynamic Symbolic Execution* (DSE) avec symbolisation des arguments
- 4 Re-concretisation des symboles (non-pointeur) pour éviter l'explosion combinatoire
- 5 Traduction en *C-like* des accès mémoires
- 6 Génération du test Python

Exemple : OpenSSL SHA1_Update

- Binaire : sha1test
- Header : openssl/sha.h

```

1  typedef struct SHAstate_st {
2      unsigned int h0, h1, h2, h3, h4;
3      unsigned int Nl, Nh;
4      unsigned int data[16];
5      unsigned int num;
6  } SHA_CTX;
7
8  typedef int size_t;
9
10 int SHA1_Update(SHA_CTX *c, const void* data, size_t len);

```

Méthode

- Trace d'exécution

Début de fonction

RAX: 8000000000000000

RBX: 1

RCX: 7ffd9b689a40

...

Lecture en 7ffd9b689a40, 32 bits, valeur 0x12345678

...

Fin de fonction

Méthode

- Rejeu en *Dynamic Symbolic Execution* (DSE) avec symbolisation des arguments

Lecture en `@[arg0_c + 65]`, valeur `0x62`

...

Écriture en `@[arg0_c]`, valeur `0x1`

...

Méthode

- Traduction en *C-like* des accès mémoires

Lecture en `arg0_c->data[1]`, valeur `0x62`

...

Écriture en `arg0_c->Nh`, valeur `0x1`

...

Méthode

■ Génération du test Python

```

1 class TestSHA1_Update(TestHeader):
2     '''This is an auto-generated class, using the Sibyl learn module'''
3     func = 'SHA1_Update'
4
5     def init1(self):
6         # arg0_c
7         base0_ptr_size = self.field_addr("arg0_c", "(arg0_c)->num") + \
8             self.sizeof("(arg0_c)->num")
9         base0_ptr = self._alloc_mem(base0_ptr_size, read=True, write=True)
10        ...
11        # (arg1_data)[1]
12        input0_ptr = base1_ptr + self.field_addr("arg1_data", "(arg1_data)[1]")
13        ...
14        self._add_arg(0, base0_ptr) # arg0_c
15        ...
16        # *(arg1_data) = 0x61
17        self._write_mem(base1_ptr, self.pack(0x61, self.sizeof("*(arg1_data)")))
18        # (arg1_data)[1] = 0x62
19        self._write_mem(input0_ptr, self.pack(0x62, self.sizeof("(arg1_data)[1]")))
20        ...

```

Méthode

■ Génération du test Python

```

1  def check1(self) :
2      return all((
3          # Check output value
4          self._get_result() == 0x1,
5          # (arg0_c)->Nh == 0x0
6          self._ensure_mem(self.input4_ptr, self.pack(0x0, self.sizeof("(arg0_c)->Nh"))),
7          # (arg0_c)->Nl == 0x18
8          self._ensure_mem(self.input1_ptr, self.pack(0x18, self.sizeof("(arg0_c)->Nl"))),
9          # *((arg0_c)->data) == 0x636261 (without considering 0x3 offset(s))
10         self._ensure_mem_sparse(self.output0_ptr,
11                                 self.pack(0x636261,
12                                             self.sizeof("*((arg0_c)->data)")),
13                                 [0x3]),
14         # (arg0_c)->num == 0x3
15         self._ensure_mem(self.input3_ptr, self.pack(0x3, self.sizeof("(arg0_c)->num"))),
16     ))
17     ...
18 TestSetTest(init1, check1) & TestSetTest(init2, check2) & ...
  
```

Démonstration

- Apprentissage de SoftFp^a (bibliothèque de softfloat)
- Recherche dans un binaire obscurci (data-flow + control-flow) utilisant cette bibliothèque

a. SoftFp, release du 20-12-2016 : <https://bellard.org/softfp/softfp-2016-12-20.tar.gz>

En théorie

Sibyl

- est *complet* : signature connue + fonction présente \rightarrow fonction trouvée
- n'est pas *correct* : possibilité de faux positif^a

a. problème impossible dans le cas général sans parcourir l'espace des entrées

Limite en pratique

- Erreur du moteur d'émulation
- Erreur dans l'écriture du test
- Mauvaise ABI

Limites en pratique

- Appel de fonction externes / malloc
 - PoC, en cours d'implémentation !

Limites en pratique

- Appel de fonction externes / malloc
 - PoC, en cours d'implémentation !
- Fonction non apprenable
 - arithmétique de pointeur ($ptr < 0x7ffffff$)
 - $((ptr \oplus 0x7ffffff) \& (ptr \oplus (ptr + (-0x7ffffff)))) \oplus (ptr + (-0x7ffffff)) \oplus ptr \oplus 0x7ffffff).msb$

Limites en pratique

- Appel de fonction externes / malloc
 - PoC, en cours d'implémentation !
- Fonction non apprenable
 - arithmétique de pointeur ($ptr < 0x7ffffff$)
 - $((ptr \oplus 0x7ffffff) \& (ptr \oplus (ptr + (-0x7ffffff)))) \oplus (ptr + (-0x7ffffff)) \oplus ptr \oplus 0x7ffffff).msb$
- Il faut que ce soit une fonction connue de Sibyl !

Futur

- Plus de test !
 - Test sur de nouvelles bibliothèques
 - Complétion des tests d'apprentissage (mutants)

Futur

- Plus de test !
 - Test sur de nouvelles bibliothèques
 - Complétion des tests d'apprentissage (mutants)
- Plus robuste !
 - Amélioration des capacités d'apprentissage
 - Gestion des appels aux fonctions externes

Futur

- Plus de test !
 - Test sur de nouvelles bibliothèques
 - Complétion des tests d'apprentissage (mutants)

- Plus robuste !
 - Amélioration des capacités d'apprentissage
 - Gestion des appels aux fonctions externes

- Plus de fonctionnalités !
 - Détection automatique de l'ABI
 - Détection des adresses de fonctions
 - Mode « pourquoi j'ai tort ? »

Sibyl

- github.com/cea-sec/Sibyl
- [README.MD#Basic usage](#)
- [/doc](#)
- Issue + PR !

Commissariat à l'énergie atomique et aux énergies alternatives

CEA/DAM

Centre de Bruyères-le-Châtel | 91297 Arpajon Cedex

T. +33 (0)1 69 26 40 00 | F. +33 (0)1 69 26 40 00

Établissement public à caractère industriel et commercial

RCS Paris B 775 685 019