

# Finding library subroutines in stripped statically-linked binaries

findmagic

Katharina Bogad  
Technische Universität München  
Computer Science Department

SS 2015  
January 18, 2017

- ▶ Computer Science student
- ▶ Member of the H4x0rPsch0rr CTF-Team and CTF-Player for fun (and sometimes profit)
- ▶ Interested in reverse engineering for long time
- ▶ Hates QR-Codes

Who of you has...

- ▶ basic knowledge of graph theory?

Who of you has...

- ▶ basic knowledge of graph theory?
- ▶ reverse engineered a statically linked binary at least once?

- ▶ Traditional pattern-matching: exact library needed for decent results
- ▶ Works reasonably well in homogenous environments like MSVCRT
- ▶ Open source libraries?
- ▶ Embedded devices?

So, what are we doing if we cannot have symbols?

- ▶ Looking at the arguments?
- ▶ Looking at suspicious constants?  
Think of `0x8080808080` for `strlen(3)`

Let's automate this!

However, there are caveats:

- ▶ Finding arguments is not a trivial task.
- ▶ What makes a constant suspicious?

However, there are caveats:

- Finding arguments is not a trivial task.
- What makes a constant suspicious?

But automating gives new perspectives:  
Comparing callgraphs!



- ▶ Program is a set of attributed graphs  $G = (N, B)$
- ▶ Nodes  $N$  are functions
- ▶ Branches  $B$  are calls between functions

We need:

- ▶ A string definition

$$\begin{aligned} & (\forall (i, c) \in str : c \geq 0x20 \wedge c \leq 0xDF \\ & \quad \vee c = 0x0A \vee c = 0x0D \vee c = 0x09 \vee c = 0x00) \\ & \wedge |str| > 1 \\ & \wedge (\forall (i, c) \in str \mid i = \max(i, str) : c = 0x00) \\ & \wedge (\forall (i, c) \in str \mid i \neq \max(i, str) : c \neq 0x00) \end{aligned} \tag{1}$$

We need:

- ▶ A string definition

$$\begin{aligned} & (\forall (i, c) \in str : c \geq 0x20 \wedge c \leq 0xDF \\ & \quad \vee c = 0x0A \vee c = 0x0D \vee c = 0x09 \vee c = 0x00) \\ & \wedge |str| > 1 \\ & \wedge (\forall (i, c) \in str \mid i = \max(i, str) : c = 0x00) \\ & \wedge (\forall (i, c) \in str \mid i \neq \max(i, str) : c \neq 0x00) \end{aligned} \tag{1}$$

Printable characters from extended ASCII ...

We need:

- ▶ A string definition

$$\begin{aligned} & (\forall (i, c) \in str : c \geq 0x20 \wedge c \leq 0xDF \\ & \quad \vee c = 0x0A \vee c = 0x0D \vee c = 0x09 \vee c = 0x00) \\ & \wedge |str| > 1 \\ & \wedge (\forall (i, c) \in str \mid i = \max(i, str) : c = 0x00) \\ & \wedge (\forall (i, c) \in str \mid i \neq \max(i, str) : c \neq 0x00) \end{aligned} \tag{1}$$

... and `\n`, `\r`, `\t` and `0x00` ...

We need:

- ▶ A string definition

$$\begin{aligned} & (\forall (i, c) \in str : c \geq 0x20 \wedge c \leq 0xDF \\ & \quad \vee c = 0x0A \vee c = 0x0D \vee c = 0x09 \vee c = 0x00) \\ & \wedge |str| > 1 \\ & \wedge (\forall (i, c) \in str \mid i = \max(i, str) : c = 0x00) \\ & \wedge (\forall (i, c) \in str \mid i \neq \max(i, str) : c \neq 0x00) \end{aligned} \tag{1}$$

... with a minimum length of 2 ...

We need:

- ▶ A string definition

$$\begin{aligned} & (\forall (i, c) \in str : c \geq 0x20 \wedge c \leq 0xDF \\ & \quad \vee c = 0x0A \vee c = 0x0D \vee c = 0x09 \vee c = 0x00) \\ & \wedge |str| > 1 \\ & \wedge (\forall (i, c) \in str \mid i = \max(i, str) : c = 0x00) \\ & \wedge (\forall (i, c) \in str \mid i \neq \max(i, str) : c \neq 0x00) \end{aligned} \tag{1}$$

... where the last character is 0x00 and no other character is 0x00.

We need:

- ▶ A node definition

$$N = (n, s, C, S, I)$$

- ▶  $n$ : Function name
- ▶  $s$ : Function address
- ▶  $C$ : Multiset of constant values
- ▶  $S$ : Multiset of cross-referenced strings
- ▶  $I$ : Ordered multiset of the machine instructions

Objective: Generate a bijective mapping  $M = N_1 \rightarrow N_2$

- ▶  $N_1$ : known library function
- ▶  $N_2$ : function inside the target library



- 1 Acquire target library with debug symbols

- 1 Acquire target library with debug symbols
- 2 Build the graphs for it

- 1 Acquire target library with debug symbols
- 2 Build the graphs for it
- 3 Build graphs for the binary we analyse

- 1 Acquire target library with debug symbols
- 2 Build the graphs for it
- 3 Build graphs for the binary we analyse
- 4 Match them

Do we need exactly the same binary used for linking?

- ▶ Short answer: no.

Do we need exactly the same binary used for linking?

- ▶ Short answer: no.
- ▶ Long answer: it depends.

- ▶ A reasonably close version is enough
- ▶ Watch out for compiler flags
- ▶ Also problematic: `assert()`

## Caution: real world example

```
2391 assert((unsigned long) (old_size) < (unsigned long) (nb + MINSIZE));
```

with relocation:

```
1 (unsigned long) (old_size) < (unsigned
   long) (
2   nb + (unsigned long)(
3     (((__builtin_offsetof (struct
         malloc_chunk, fd_nextsize)) +
4       (
5         (2 * (sizeof(size_t))) - 1
6       ))
7     & ~(
8       (2 * (sizeof(size_t))) - 1
9     ))))
```

No code, but debug strings vary!

without relocation:

```
1 (unsigned long) (old_size) < (unsigned
   long) (
2   nb + (unsigned long)(
3     (((__builtin_offsetof (struct
         malloc_chunk, fd_nextsize)) +
4       ((2 * (sizeof(size_t)) <
5         __alignof__ (long double) ?
6         __alignof__ (long double) :
7         2 * (sizeof(size_t))
8       ) - 1))
9     & ~(
10      (2 * (sizeof(size_t)) <
11        __alignof__ (long double) ?
12        __alignof__ (long double) :
13        2 * (sizeof(size_t))
14      ) - 1
15    ))))
```



- 1 Iterate over subroutines
- 2 Iterate over the instructions of these subroutines
- 3 If something interesting is found, add it to the corresponding list<sup>1</sup>

---

<sup>1</sup>See the paper for a marvellous formal definitions for this

- ▶ `call` instructions add a new branch to the functions callgraph
- ▶ Additionally for Intel x86\_64 architecture:
- ▶ Only if it's a near call - opcode `0xE8`
- ▶ This ensures we're in the same section
- ▶ Other architectures may need different conditions!

- ▶ Look for something that loads a pointer (x86\_64: `lea, mov`)
- ▶ Check if it's a string by our definition
- ▶ If so, add it to the Strings of the current function

- ▶ We don't want to add pointer arithmetic as constants
- ▶ Interesting constants are often bitmasks
- ▶ Thus, we limit ourselves to the immediates of `and`, `or`, `xor` and `mov`
- ▶ Optionally, we may exclude further by doing value checking on the constant

Isomorphism:

- ▶ Ancient greek: *isos* = equal and *morphe* = shape
- ▶ Mathematical way to compare the structure of objects

Choosing the right algorithm:

- ▶ Ullmann's algorithm
- ▶ Nauty (**no automorphism, yes?**)
- ▶ VF2

Choosing the right algorithm:

- ▶ Ullmann's algorithm
- ▶ Nauty (no automorphism, yes?)
- ▶ VF2

Choosing the right algorithm:

- ▶ Ullmann's algorithm
- ▶ Nauty (no automorphism, yes?)
- ▶ VF2



Choosing the right algorithm:

- ▶ Ullmann's algorithm
- ▶ Nauty (**no automorphism, yes?**)
- ▶ VF2

Choosing the right algorithm:

- ▶ Callgraphs cannot be considered randomly connected
- ▶ Some functions imply calls to other functions
- ▶ `malloc()` & `free()`, `accept()` & `close()`, ...
- ▶ VF2 is very fast in this situation
- ▶ Also, VF2 can check semantic equality of the nodes in the same step

- ▶  $G_1 = (N_1, B_1), G_2 = (N_2, B_2)$
- ▶ Mapping  $M \subset N_1 \times N_2$
- ▶  $M$  must be a bijective function
- ▶  $M$  must not alter the branch structure

- ▶ State Space Representation (SSR)  $s$
- ▶ Essentially a set of tuples  $(n_1, n_2)$
- ▶  $M(s)$  denotes a partial mapping
- ▶ Two subgraphs  $G_1(s)$  and  $G_2(s)$  can be derived, containing only the nodes in the matching and the branches connecting them
- ▶ Same for  $M_1(s), M_2(s), B_1(s), B_2(s)$

- ▶ Transition from state  $s$  to  $s'$ :  $s' = s \cup \{(n, m)\}$
- ▶ But: only a small set of these states are *consistent*
- ▶ We introduce *k-lookahead rules* to conclude whether a consistent state can be reached after  $k$  steps
- ▶ These rules will be called *feasibility rules*

Feasibility function:  $F(s, n, m) = F_{syn}(s, n, m) \wedge F_{sem}(s, n, m)$

- ▶  $F_{syn} \rightarrow$  **syntactic** feasibility
- ▶  $F_{sem} \rightarrow$  **semantic** feasibility

- ▶ Initial state is empty, i.e.  $M(s_0) = \emptyset$
- ▶ In each step, compute  $P(S)$ , the node pairs of candidates to be added
- ▶  $T_n^{in} \rightarrow$  nodes with branches ending into  $G_n(s)$
- ▶  $T_n^{out} \rightarrow$  nodes with branches starting from  $G_n(s)$
- ▶  $P(s) = \{(n, m) | n \in T_1^{out}, m \in T_2^{out}\}$  if no  $T_n^{out}$  is empty,  $T_n^{in}$  otherwise
- ▶ If  $P(s)$  is still empty, backtrack until a state  $s$  is reached with  $P(s)$  containing not examined node pairs

Check predecessors of current node:

$$\begin{aligned} R_{pred}(s, n, m) &\iff \\ &(\forall n' \in M_1(s) \cap \text{Pred}(G_1, n) \exists m' \in \text{Pred}(G_2, m) \mid (n', m') \in M(s)) \wedge \\ &(\forall m' \in M_2(s) \cap \text{Pred}(G_2, m) \exists n' \in \text{Pred}(G_1, n) \mid (n', m') \in M(s)) \end{aligned}$$



Check successors of current node:

$$\begin{aligned} R_{succ}(s, n, m) &\iff \\ &(\forall n' \in M_1(s) \cap \text{Succ}(G_1, n) \exists m' \in \text{Succ}(G_2, m) \mid (n', m') \in M(s)) \wedge \\ &(\forall m' \in M_2(s) \cap \text{Succ}(G_2, m) \exists n' \in \text{Succ}(G_1, n) \mid (n', m') \in M(s)) \end{aligned}$$

1-lookahead:

$$R_{in}(s, n, m) \iff$$

$$(|\text{Succ}(G1, n) \cap T_1^{in}(s)| = |\text{Succ}(G2, m) \cap T_2^{in}(s)|) \wedge$$

$$(|\text{Pred}(G1, n) \cap T_1^{in}(s)| = |\text{Pred}(G2, m) \cap T_2^{in}(s)|)$$

$$R_{out}(s, n, m) \iff$$

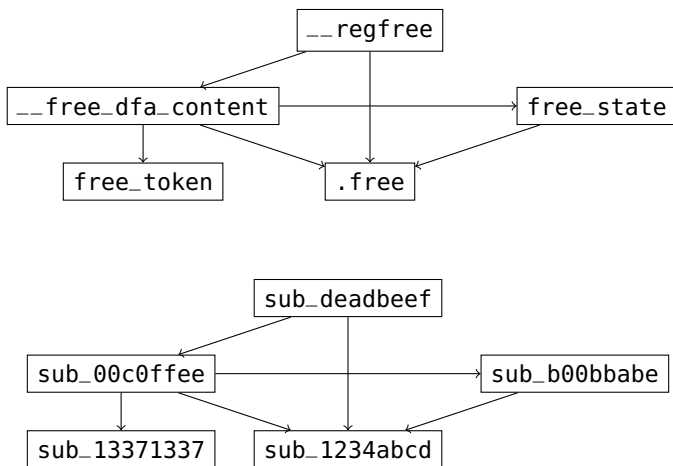
$$(|\text{Succ}(G1, n) \cap T_1^{out}(s)| = |\text{Succ}(G2, m) \cap T_2^{out}(s)|) \wedge$$

$$(|\text{Pred}(G1, n) \cap T_1^{out}(s)| = |\text{Pred}(G2, m) \cap T_2^{out}(s)|)$$

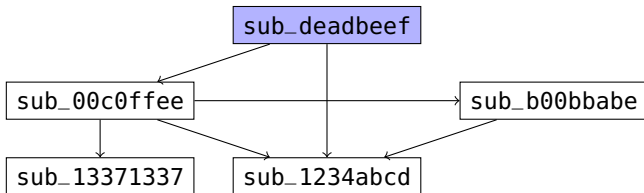
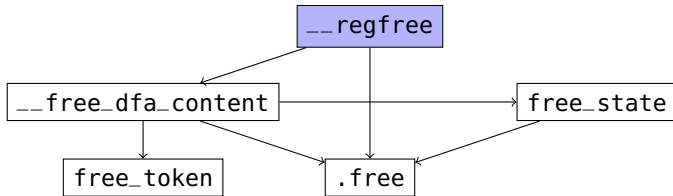
2-lookahead:

$$\begin{aligned} R_{new}(s, n, m) &\iff \\ &(|\tilde{N}_1(s) \cap \text{Pred}(G_1, n)| = |\tilde{N}_2(s) \cap \text{Pred}(G_2, m)|) \wedge \\ &(|\tilde{N}_1(s) \cap \text{Succ}(G_1, n)| = |\tilde{N}_2(s) \cap \text{Succ}(G_2, m)|) \end{aligned}$$

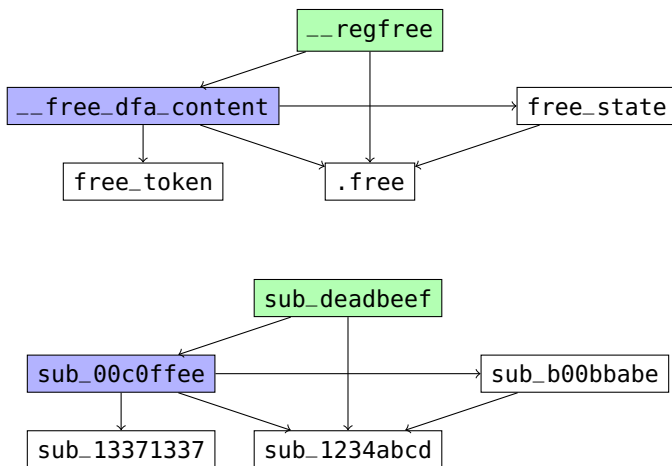
Example:



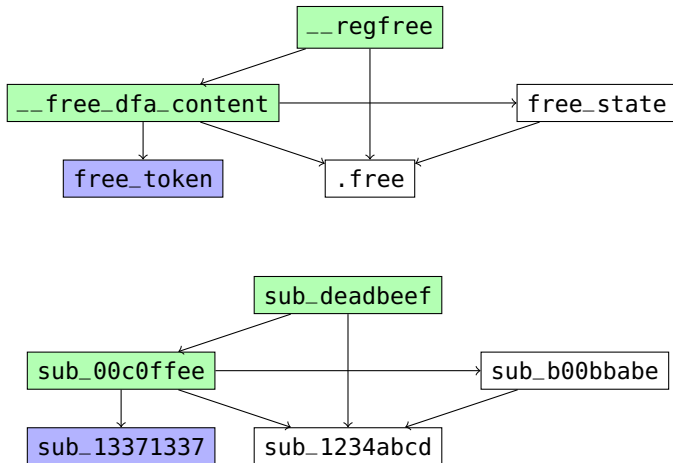
Example:



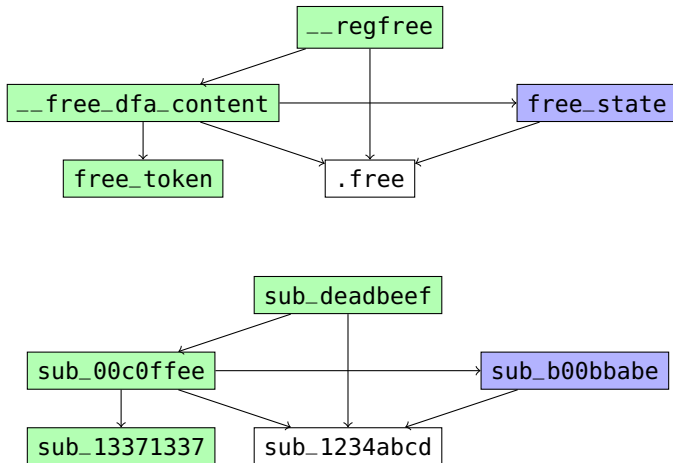
Example:



Example:

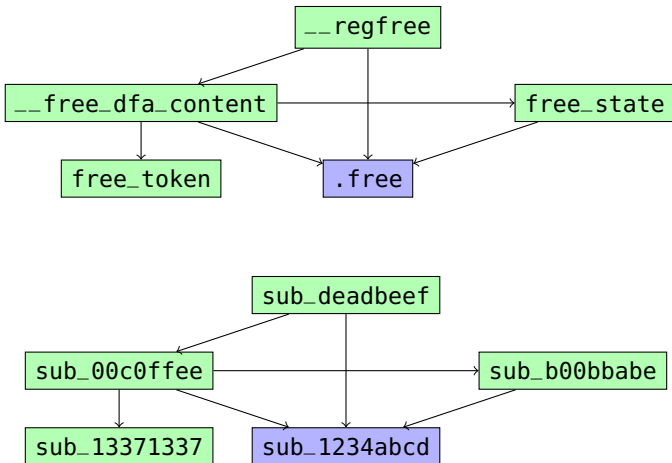


Example:

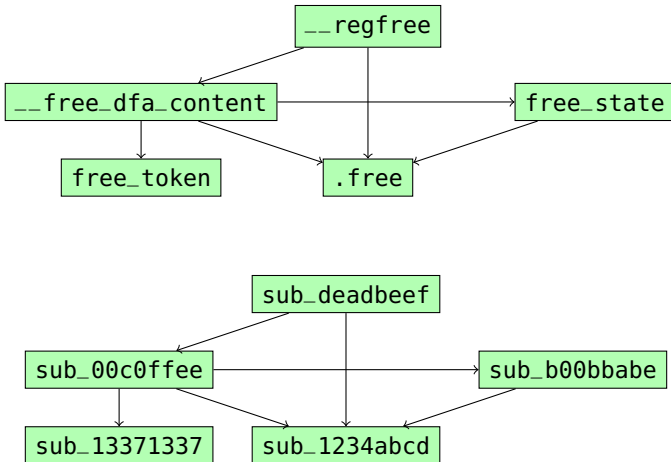




Example:



Example:



Semantic feasibility: We define a compatibility relation  $\approx$

$$\begin{aligned}n \approx m &\iff \\&(\forall c \in C_n \exists c' \in C_m | c = c') \wedge \\&(\forall c \in C_m \exists c' \in C_n | c = c') \wedge \\&(\forall s \in S_n \exists s' \in S_m | s = s') \wedge \\&(\forall s \in S_m \exists s' \in S_n | s = s')\end{aligned}$$

This yields the final rule:

$$F_{sem}(s, n, m) \iff \\ n \approx m \wedge \forall (n', m') \in M(s), (n, n') \in B_1 \Rightarrow (n, n') \approx (m, m') \\ \wedge \forall (n', m') \in M(s), (n', n) \in B_1 \Rightarrow (n', n) \approx (m', m)$$

- ▶ Matching is done in brute-force manner
- ▶ Multiple sets:
  - ▶ Functions that can be exactly identified
  - ▶ Functions that have multiple, possible matches
  - ▶ Functions that cannot be found via matching (no strings, no constants and no function calls - `IO_default_sync`)

- ▶ Test implementation was created
- ▶ Free as in Speech (GPLv3 or Later)
- ▶ Grab it from github:  
<https://github.com/masterofjellyfish/findmagic>
- ▶ Disclaimer: You need .NET Framework or Mono
- ▶ Supports only x86\_64 for now
- ▶ Major code cleanup and more architectures (ARM, MIPS) are planned

Exact matches:

definitions	linked against	find- magic	FLIRT (IDA)
glibc 2.21 (arch linux)	glibc 2.21 (arch linux)	376	233
glibc 2.21 (arch linux)	glibc 2.13 (debian- wheezy)	105	72

- ▶ Algorithm can also provide hints
- ▶ For example: `strcpy_sse2`, `strcpy_sse3`
- ▶ Same constants, same callgraph
- ▶ Indistinguishable by the algorithm, but they do the same job
- ▶ Helpful for manual reversing!



Recovery fails if multiple matching possibilities and function is not part of unique call graph. Example:

```
1 .CapstoneX86Detail
2   push rbp
3   mov rbp, rsp
4   mov rax, rdi
5   add rax, 0x30
6   leave
7   retn
```

Pseudocode:

```
1 cs_x86* CapstoneX86Detail(cs_detail *detail) {
2   return &detail->x86;
3 }
```

Thanks!