# Binary Analysis - 2

**REVERSE ENGINEERING**

deti **universidade de aveiro**
departamento de eletrónica,
telecomunicações e informática

**João Paulo Barraca**

# Binary Analysis Process

- Up to now we know how ELF files are structure, but the question remains: how do we analyse ELF files?
  - Or any other binary

- A possible flow can be:
  - File analysis (file, nm, ldd, content visualization, foremost, binwalk)

  - Static Analysis (disassemblers and decompilers

  - Behavioral Analysis (strace, LD_PRELOAD)

  - Dynamic Analysis (debuggers)

# Identifying a file

- Files should be seen as containers (this includes ELF files)
  - May have the expected content type
    - But it may have an unexpected behavior (e.g. bug or malware)
  - May have unexpected, additional content (e.g. polyglots)
    - More common in DRM schemes and malware in order to hide binary blobs

- Files should not be trusted
  - Both the expected and additional content may be malicious
  - Static analysis is safe (as long as nothing is executed)
  - Dynamic analysis is not safe. Sandboxes and VMs must be used

# Questions to answer

- What type of file we have?
  - Are there hidden contents?

- What is the architecture?

- Is it 64/32 or ARM7/ARM9/ARM9E/ARM10?

- Where is the starting address?

- What the main function does?

- What will the program will actually do?

# Questions to answer

## Some basic tools go a long way

- **`file`**: (try to identify) the type of file
  - Only applies to a top container. File is not able to look into enclosed binary blobs
  - Alternatives that complement **`file`** are **`binwalk`** and **`foremost`**

- **`xxd`**: hexdump the file, allowing to rapidly detect patterns
  - less also helps to hold the content in the terminal

- **`strings`**: prints null terminated sequence chars
  - By default, with more than 4 characters (-n setting)

- **`ldd`**: print shared object dependencies
  - Libraries registered in the ELF that are required (typically for dynamically relocate symbols)

- **`nm`**: dumps symbols from `.symtab` (or `.dyntab` with –D)

# Disassembler basics with ghidra

- ghidra is a open source tool developed by NSA and released to the public doing Disassembly and Static Analysis
  - Development branch has support for Dynamic Analysis (should be released "soon")

- Works on Windows, Linux and macos
  - Java based

- Not the most important tool (IDA is), but is gaining a huge traction
  - It's free and very powerful with a huge number of platforms and a **fine decompiler**

CodeBrowser: project:/authenticator

File  Edit  Analysis  Graph  Navigation  Search  Select  Tools  Window  Help

**Top menu** and tools for quick access.

Program Trees

authenticator
 .bss
 .data
 .got
 .dynamic
 .fini_array
 .init_array
 .eh_frame
 .eh_frame_hdr
 .rodata

Program Tree ×

Symbol Tree

 Imports
 Exports
 Functions
  __cxa_finalize
  __cxa_finalize
  __do_global_dtors_aux
  __gmon_start__
  __libc_csu_fini
  __libc_csu_init
  __libc_start_main
  __stack_chk_fail
  __stack_chk_fail
  _fini

Filter:

Data Type Manager

 Data Types
  BuiltInTypes
  authenticator
  generic_clib
  generic_clib_64
  jni_all
  libCPython
  lua

Filter:

Listing: authenticator

```
                                         //
                                         // segment_2.1
                                         // Loadable segment  [0x0 - 0xelf] (disabled execute bit)
                                         // ram:00100000-ram:00100237
                                         //
                           01 00 00 ...
        00100000 7f            db        7Fh                 e_ident_magi...
        00100001 45 4c 46      ds        "ELF"               e_ident_magi...
        00100004 02            db        2h                  e_ident_class
        00100005 01            db        1h                  e_ident_data
        00100006 01            db        1h                  e_ident_vers...
        00100007 00            db        0h                  e_ident_osabi
        00100008 00            db        0h                  e_ident_abiv...
        00100009 00 00 00 00 00  db[7]                       e_ident_pad
                 00 00
        00100010 03 00         dw        3h                  e_type
        00100012 3e 00         dw        3Eh                 e_machine
        00100014 01 00 00 00   ddw       1h                  e_version
        00100018 d0 07 00 00 00  dq                _start    e_entry
                 00 00
        00100020 40 00 00 00 00  dq       Elf64_Phdr_ARRAY_00100... e_phoff    =
                 00 00
        00100028 00 2b 00 00 00  dq       Elf64_Shdr_ARRAY_elfS... e_shoff
                 00 00
        00100030 00 00 00 00    ddw       0h                  e_flags
        00100034 40 00         dw        40h                 e_ehsize
        00100036 38 00         dw        38h                 e_phentsize
        00100038 09 00         dw        9h                  e_phnum
        0010003a 40 00         dw        40h                 e_shentsize
        0010003c 1d 00         dw        1Dh                 e_shnum
        0010003e 1c 00         dw        1Ch                 e_shstrndx

                      Elf64_Phdr_ARRAY_00100040     XREF[2]:    00100020(*), 0010005
        00100040 06 00 00       Elf64_Ph...                      PT_PHDR - Program
                 00 04 00
                 00 00 40 ...
```

Decompile: _start - (authenticator)

```
1
2  void _start(undefined8 param_1,undefined8 param_2,undefined8 param_3)
3
4  {
5    undefined8 in_stack_00000000;
6    undefined auStack8 [8];
7
8    __libc_start_main(main,in_stack_00000000,&stack0x00000008,__libc_csu_init,__libc_csu_fini,pa
9                      auStack8);
10   do {
11                     /* WARNING: Do nothing block with infinite loop */
12   } while( true );
13 }
14
```

Defined Strings ×   Decompile: _start ×

Function Call Trees - <No Function>

Incoming Calls
 No Function

Outgoing Calls
 No Function

Filter:        Filter:

00100000

**Executable Structure (ELF, PE...)**

All that was previously addressed can be inspected here.

Particular relevant to check content of additional sections, .got .symtab and .dynamic

**Clicking on the file name will present the header, which contains the entry point.**

**Functions** declared in the code. Not only the user defined functions, but also additional ones injected by the compiler, or required to the OS and linker.

**Exports** are functions declared in the binary file, which are available for use by another programs. That is: a name that points to some code.

**Imports** are functions used from external libraries. It includes the libraries required for the resolution of the said names.

**Data types** are methods to organize and represent memory in the scope of the program address space. While developing a program we call these as Enums or Structures.

By clicking on the low arrow we can define new structures. The interface will ask us to state a name, and then add fields with a type. The type can be a basic one (int, short, double, or a complex one)

In File -> Parse C Source we can include a C header (.h) with types directly to our program

File   Edit   Analysis   Graph   Navigation   Search   Select   Tools   Window   Help

**Program Trees**

- authenticator
  - .bss
  - .data
  - .got
  - .dynamic
  - .fini_array
  - .init_array
  - .eh_frame
  - .eh_frame_hdr
  - .rodata

Program Tree ×

**Symbol Tree**

- Imports
- Exports
- Functions
  - __cxa_finalize
  - __cxa_finalize
  - __do_global_dtors_aux
  - __gmon_start__
  - __libc_csu_fini
  - __libc_csu_init
  - __libc_start_main
  - __stack_chk_fail
  - __stack_chk_fail
  - _fini

Filter:

**Data Type Manager**

- Data Types
  - BuiltInTypes
  - authenticator
  - generic_clib
  - generic_clib_64
  - jni_all
  - libCPython
  - lua

Filter:

**Listing: authenticator**

```
                          //
                          // segment_2.1
                          // Loadable segment   [0x0 - 0xelf] (disabled execute bit)
                          // ram:00100000-ram:00100237
                          //
              assume DF = 0x0   (Default)
     00100000 7f 45 4c        Elf64_Ehdr
              46 02 01
              01 00 00 ...
     00100000 7f             db        7Fh           e_ident_magi...
     00100001 45 4c 46       ds        "ELF"         e_ident_magi...
     00100004 02             db        2h            e_ident_class
     00100005 01             db        1h            e_ident_data
     00100006 01             db        1h            e_ident_vers...
     00100007 00             db        0h            e_ident_osabi
     00100008 00             db        0h            e_ident_abiv...
     00100009 00 00 00 00 00  db[7]                  e_ident_pad
              00 00
     00100010 03 00          dw        3h            e_type
     00100012 3e 00          dw        3Eh           e_machine
     00100014 01 00 00 00    ddw       1h            e_version
     00100018 d0 07 00 00 00 dq        _start        e_entry
     00100020 40 00 00 00 00 dq        Elf64_Phdr_ARRAY_00100... e_phoff       =
              00 00 00
     00100028 00 2b 00 00 00 dq        Elf64_Shdr_ARRAY__elfS... e_shoff
              00 00 00
     00100030 00 00 00 00    ddw       0h            e_flags
     00100034 40 00          dw        40h           e_ehsize
     00100036 38 00          dw        38h           e_phentsize
     00100038 09 00          dw        9h            e_phnum
     0010003a 40 00          dw        40h           e_shentsize
     0010003c 1d 00
     0010003e 1c 00
                    Elf64
     00100040 06 00 00
              00 04 00
              00 00 40 ...
```

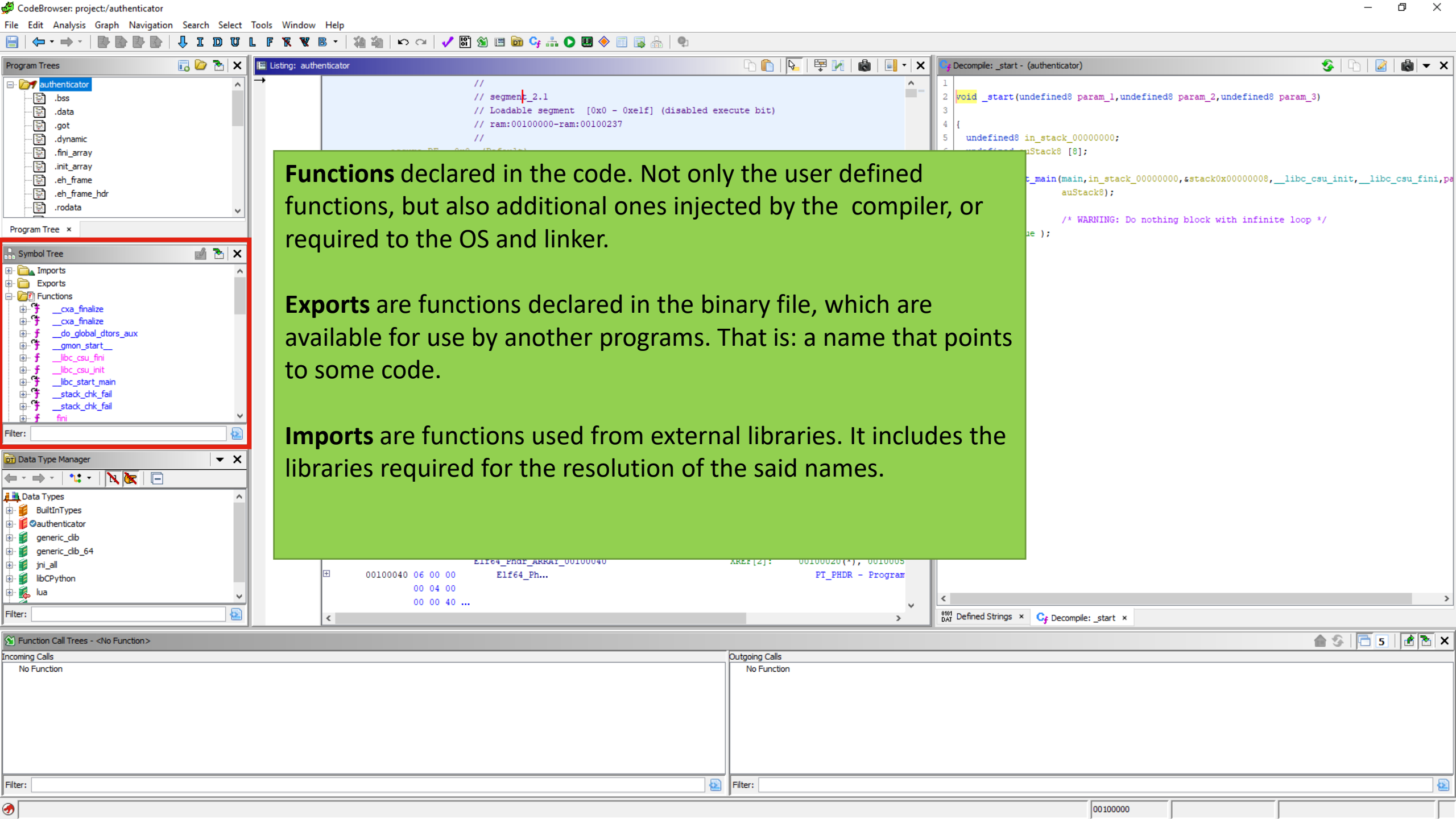**Decompile: _start - (authenticator)**

```
1
2   void _start(undefined8 param_1,undefined8 param_2,undefined8 param_3)
3
4   {
5     undefined8 in_stack_00000000;
6     undefined8 auStack8 [8];
7
8     __libc_start_main(main,in_stack_00000000,&stack0x00000008,__libc_csu_init,__libc_csu_fini,pa
9                       auStack8);
10    do {
11                /* WARNING: Do nothing block with infinite loop */
12    } while( true );
13  }
14
```
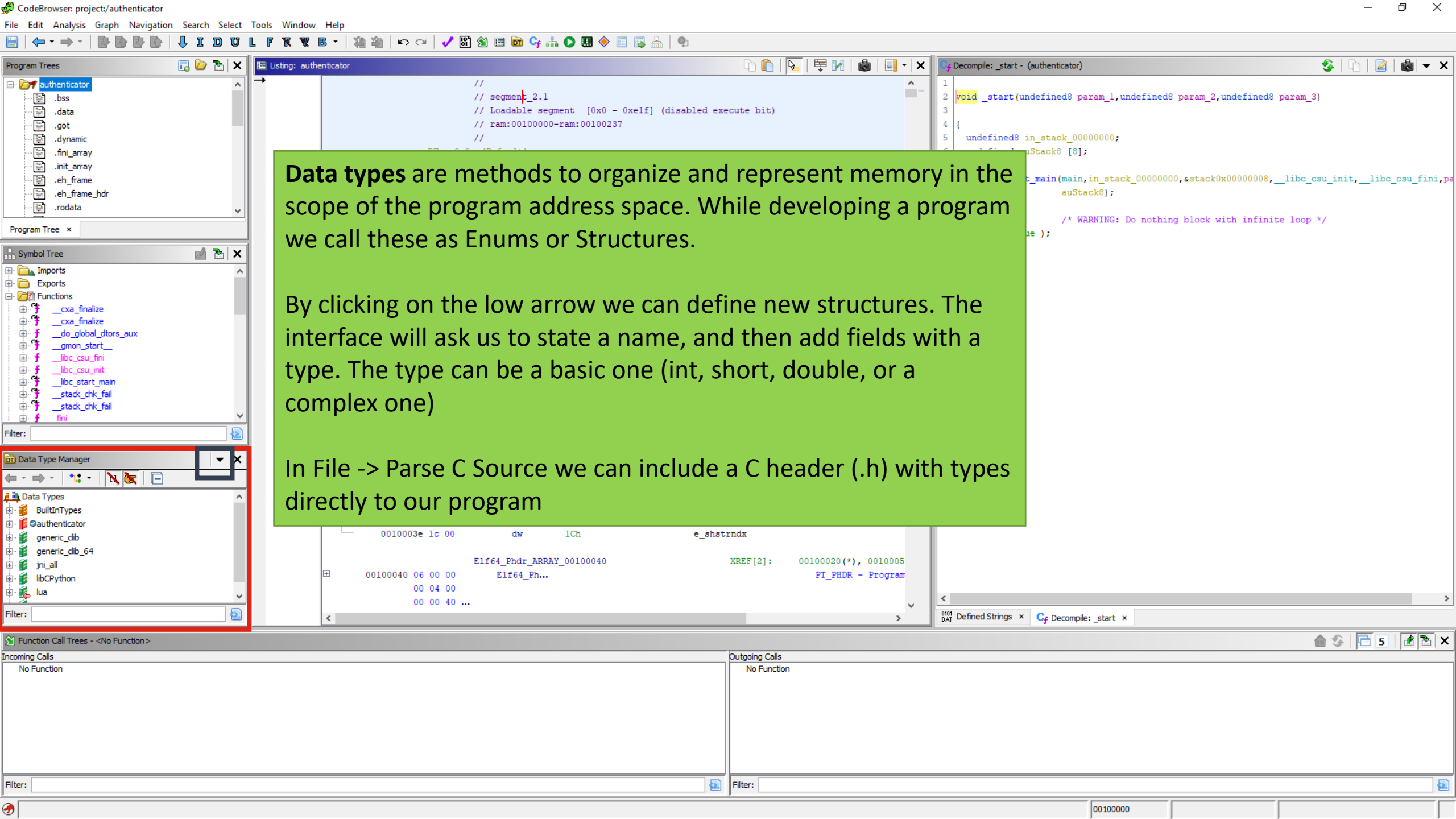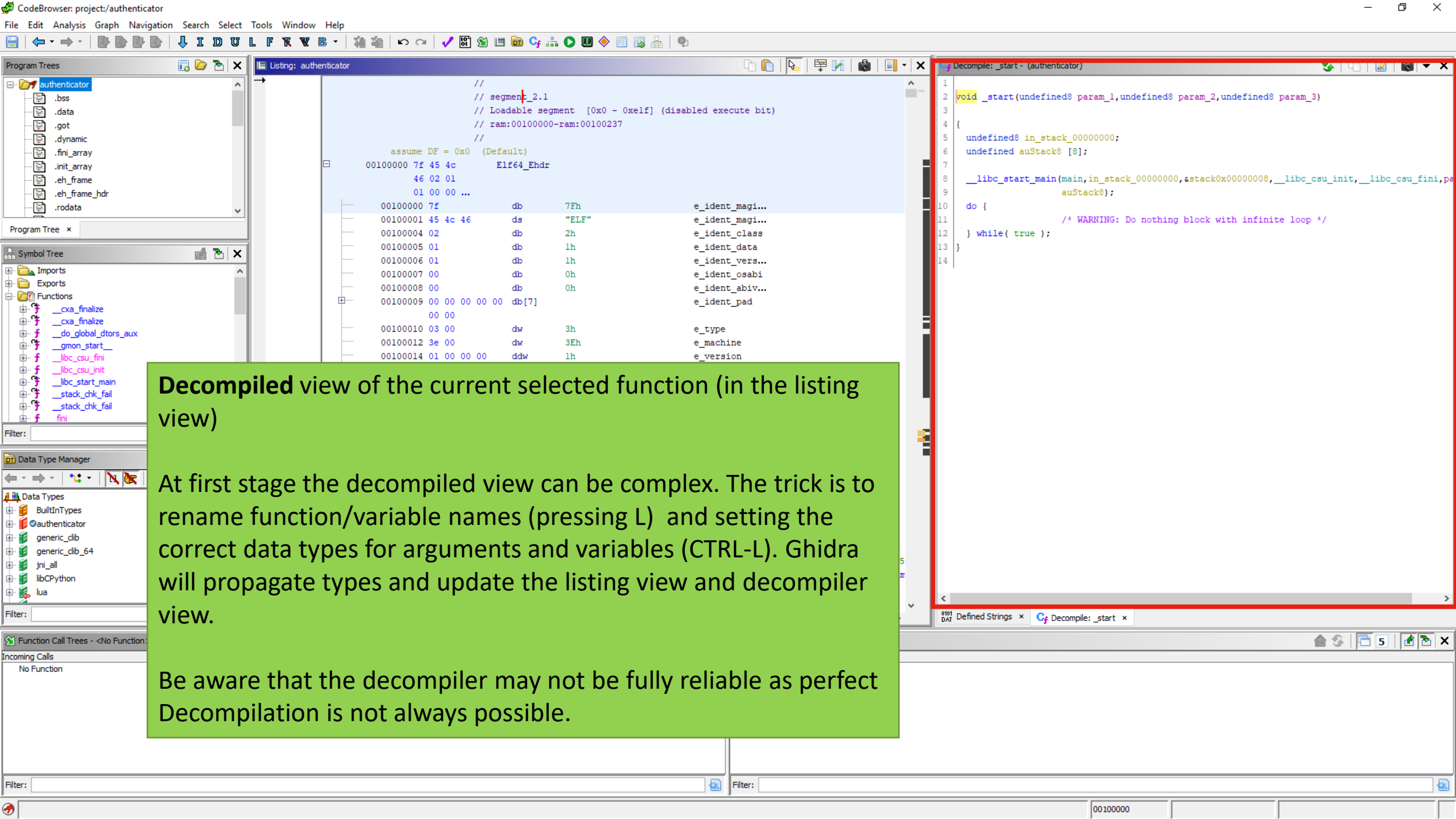
**Listing view** of the program segments.

The different bytes are displayed according to their type. Code is disassembled and the instructions are presented. For x86, the Intel format is used by default
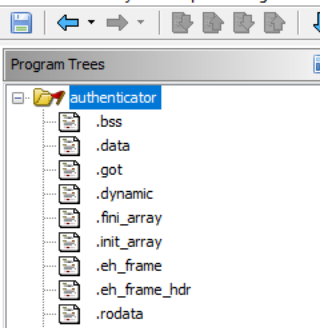
The format is something like:
Address    bytes    interpretation    comments

Incoming Calls

No Function

00100000

**Decompiled** view of the current selected function (in the listing view)

At first stage the decompiled view can be complex. The trick is to rename function/variable names (pressing L) and setting the correct data types for arguments and variables (CTRL-L). Ghidra will propagate types and update the listing view and decompiler view.
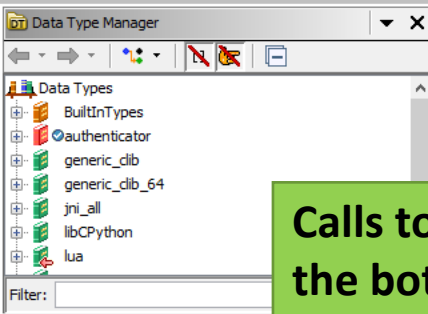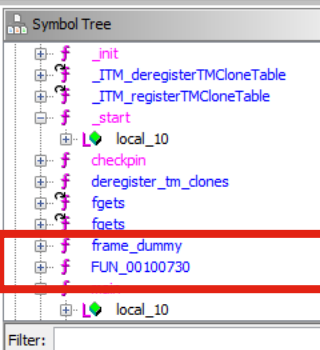
Be aware that the decompiler may not be fully reliable as perfect Decompilation is not always possible.

File   Edit   Analysis   Graph   Navigation   Search   Select   Tools   Window   Help

**Program Trees**

- authenticator
  - .bss
  - .data
  - .got
  - .dynamic
  - .fini_array
  - .init_array
  - .eh_frame
  - .eh_frame_hdr
  - .rodata

Program Tree ×

**Symbol Tree**

- _init
- _ITM_deregisterTMCloneTable
- _ITM_registerTMCloneTable
- _start
  - local_10
- checkpin
- deregister_tm_clones
- fgets
- fgets
- frame_dummy
- FUN_00100730
  - local_10

Filter:

**Data Type Manager**

- Data Types
  - BuiltInTypes
  - authenticator
  - generic_clib
  - generic_clib_64
  - jni_all
  - libCPython
  - lua

Filter:

**Listing view** presents functions with name, if the name is in the .dynamic or .symtab.

Otherwise, it will name functions as FUN_ADDRESS.

Functions can be identified by the symbols associated with an address or with assembly instructions. Functions are at address that are called and usually start by a stack ini

size_t strlen(cha

XREF[1]:     001009ae(j)

****************
*
****************

XREF[2]:     001009ec(W)
             00100ae7(R)
XREF[2]:     00100a98(*)
             00100aa9(*)
XREF[2]:     00100a40(*)
             00100a51(*)
XREF[4]:     Entry Point(*),
             _start:001007ed(*),
             00100da4(*)

main

001009db 55          PUSH    RBP
001009dc 48 89 e5    MOV     RBP,RSP
001009df 48 83 ec 50 SUB     RSP,0x50
001009e3 64 48 8b    MOV     RAX,qword ptr FS:[0x28]

**Decompile: main - (authenticator)**

```c
undefined8 main(void)

{
  int iVar1;
  long in_FS_OFFSET;
  char local_58 [32];
  char local_38 [40];
  long local_10;

  local_10 = *(long *)(in_FS_OFFSET + 0x28);
  setbuf(stdout,(char *)0x0);
  printstr(&DAT_00100bc3,0);
  printstr("Please enter your credentials to continue.\n\n",0);
  printstr("Alien ID: ",0);
  fgets(local_58,0x20,stdin);
  iVar1 = strcmp(local_58,"11337\n");
  if (iVar1 == 0) {
    printstr("Pin: ",0);
    fgets(local_38,0x20,stdin);
    iVar1 = checkpin(local_38);
```

**Listing view** presents References to memory locations, which are locations where code refers to a given memory address.

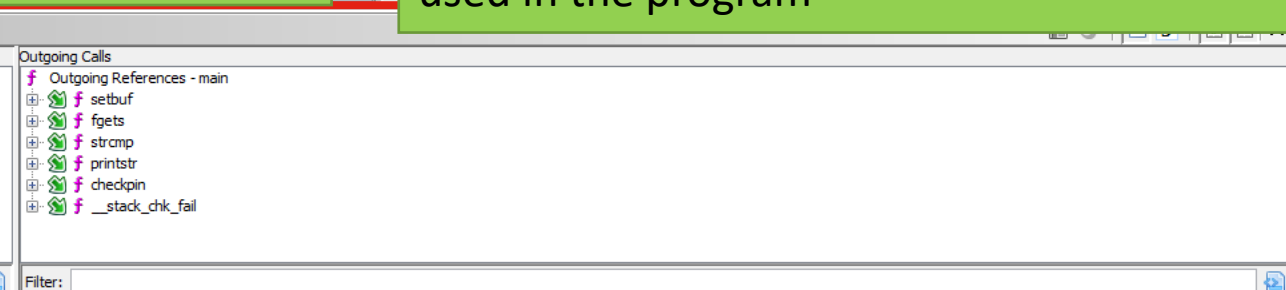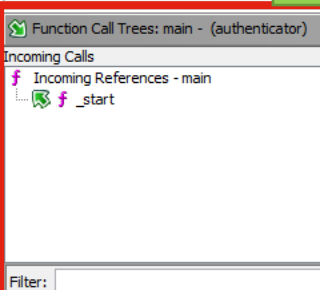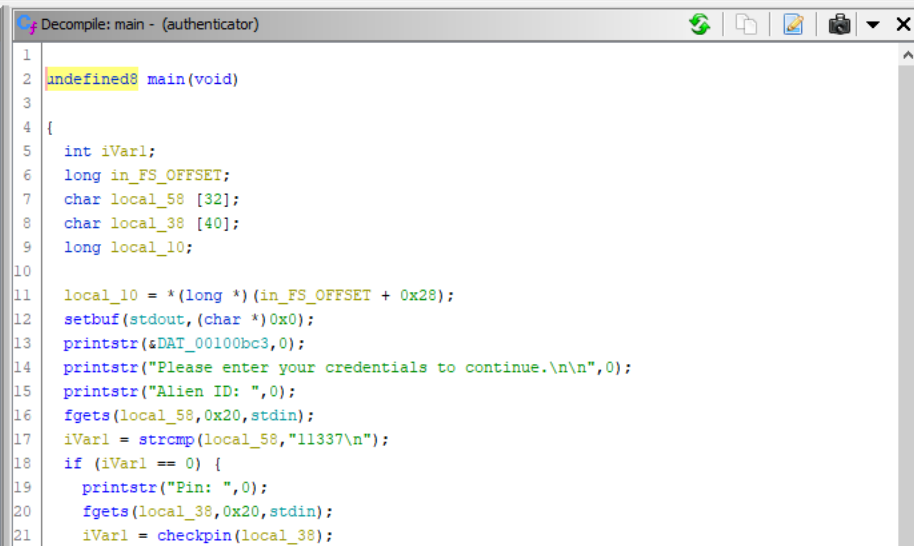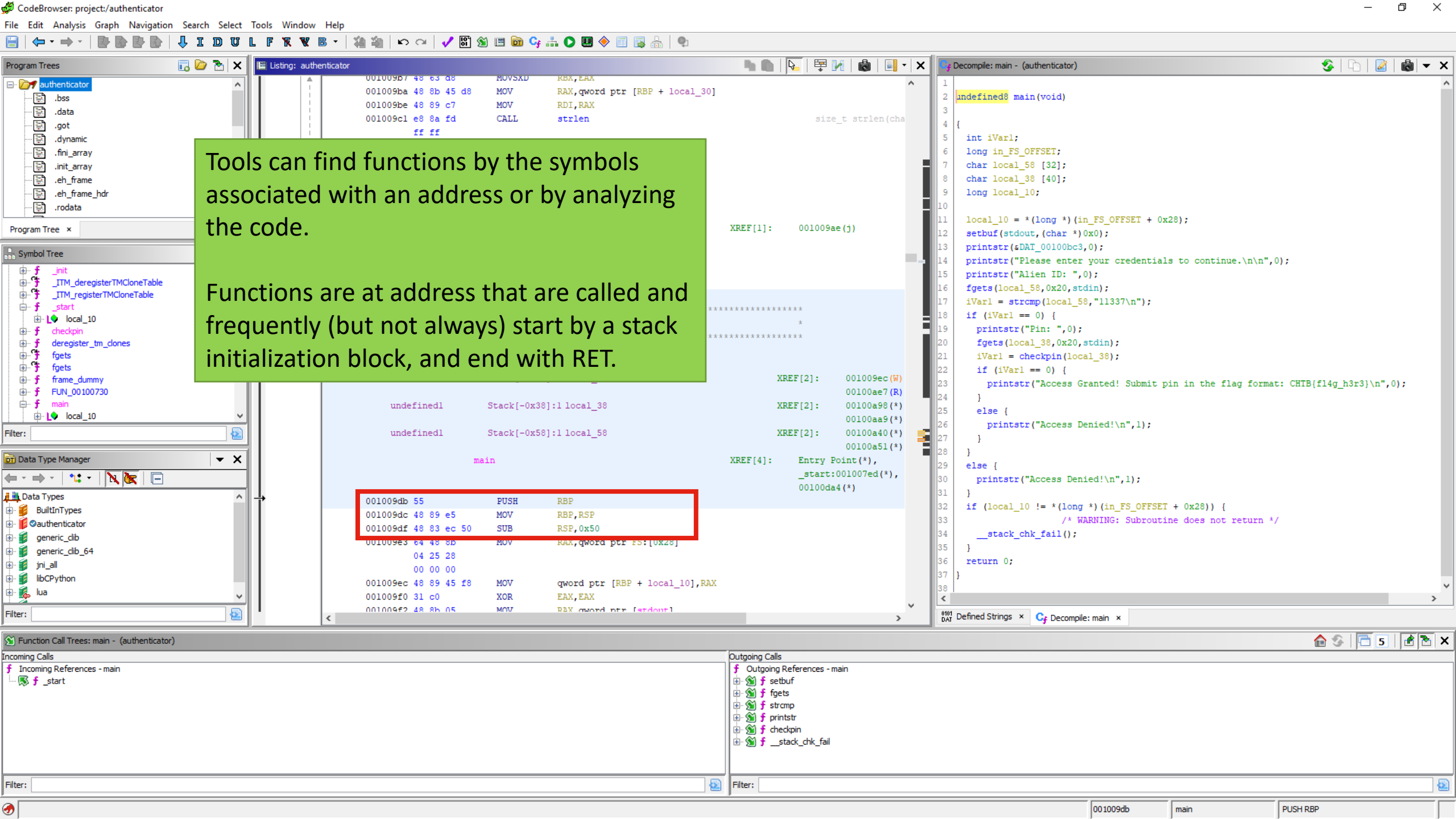May be used to identify location of arguments, function callers or data chunks used in the program

**Calls to a function or from the current function are presented in the bottom**

**Function Call Trees: main - (authenticator)**

Incoming Calls
- Incoming References - main
  - _start

Outgoing Calls
- Outgoing References - main
  - setbuf
  - fgets
  - strcmp
  - printstr
  - checkpin
  - __stack_chk_fail

Filter:

Filter:

001009db    main    PUSH RBP

Tools can find functions by the symbols associated with an address or by analyzing the code.

Functions are at address that are called and frequently (but not always) start by a stack initialization block, and end with RET.

In Menu->Window->Function Graph

A logical structure of the function is presented. This is generated by interpreting the branches that segment the function code.

**Called: Control-Flow Graphs**

In Menu->Window->Function Call Graph

A logical structure of the program is presented, starting the current function. At each node we can Show/Hide Calling functions or Called functions. We can effectively have a full representation of the program structure.

HINT: It makes much more sense with symbols or renamed functions

**Called: Call Graphs**

REVERSE ENGINEERING

# CFGs

- It is useful to think of machine code in a graph structure, called a control-flow graph

- A node in a CFG is a group of adjacent instructions called a **basic block**:
    - The only jumps into a basic block are to the first instruction
    - The only jumps out of a basic block are from the last instruction
    - I.e., a basic block always executes as a unit

- Edges between blocks represent possible jumps

# CFGs

- Basic block *a* dominates basic block *b* if every path to *b* passes through *a* first
    - strictly dominates if *a != b*


- Basic block *b* post-dominates *a* if every path through *a* also passes through *b* later

# Disassembly

- The disassembly process involves analyzing the binary, converting binary code to assembly
  - But "binary" is just a sequence of bytes, that must be mapped in the scope of a given architecture
  - Conversion depends on many factors, including compiler and flags

- Process is not perfect and may induce RE Analysts in error
  - Present instructions that actually do not exist
  - Ignore instructions that are in the binary code

- Main approaches:
  - Linear Disassembly
  - Recursive Disassembly

# Linear Disassembly

- Simplest approach towards analyzing a program: **Iterate over all code segments, disassembling the binary code as opcodes are found**

- Start at some address and follow the binary
  - Entry point or other point in the binary file
  - Entry point may not be known

- Works best with:
  - binary blobs such as from firmwares (start at the beginning)
  - objects which do not have data at the beginning
  - architecture uses variable length instructions (x86)
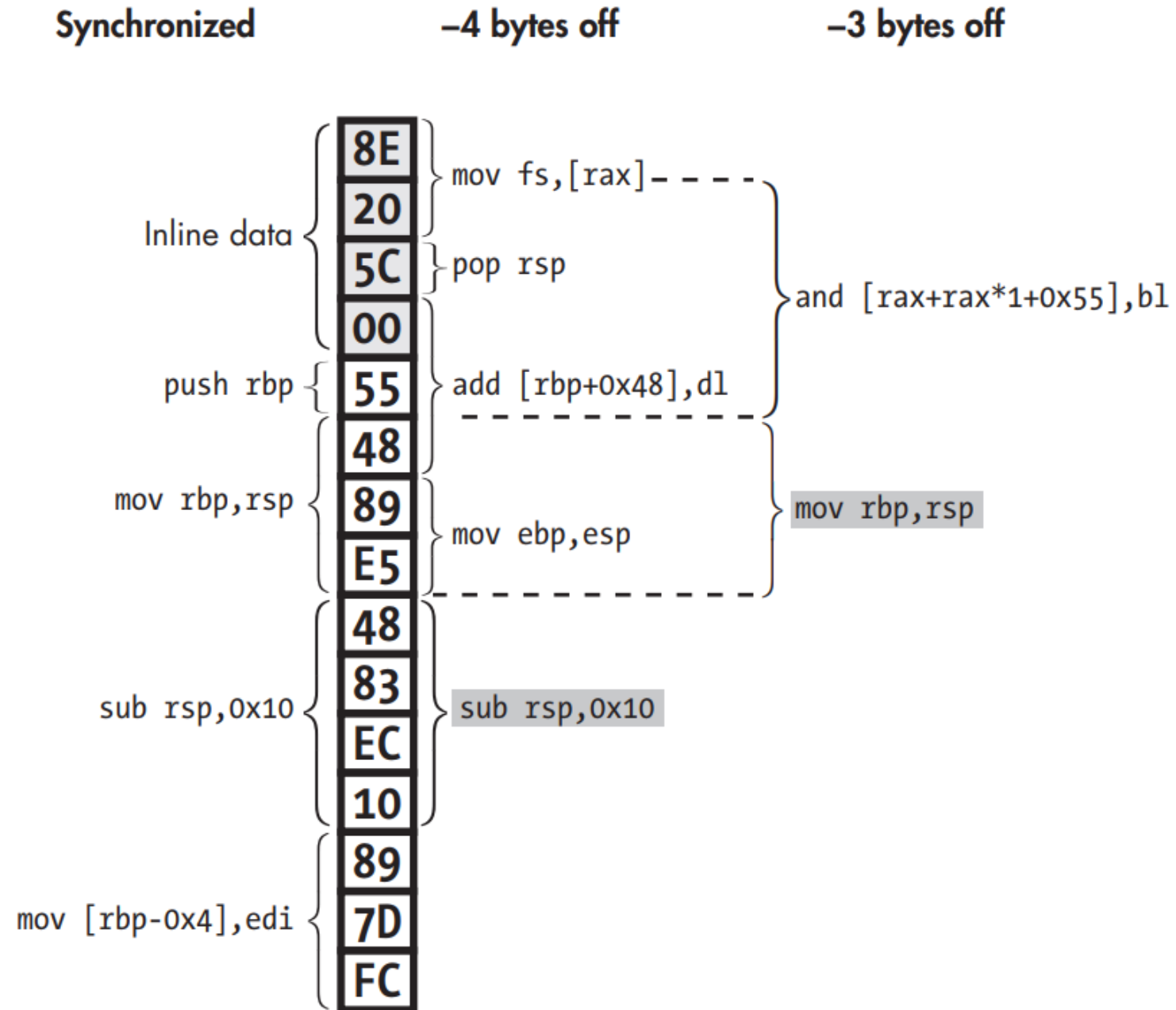
# Linear Disassembly

It is vital to define the initial address for decompiling.

An offset error will result in invalid or wrong instructions being decoded.

Linear disassembly will also try to disassemble data from the binary as if it was actual code.

Linear Disassembly is oblivious to the actual Program Flow.

With x86, because it each opcode has a variable length, the code tends to auto synchronize, but the first instructions will be missed

| Synchronized | −4 bytes off | −3 bytes off |
|---|---|---|

```
Inline data {  8E    mov fs,[rax]----
               20
               5C    pop rsp                    and [rax+rax*1+0x55],bl
               00
push rbp    {  55    add [rbp+0x48],dl
               48
mov rbp,rsp {  89    mov ebp,esp                mov rbp,rsp
               E5
               48
               83    sub rsp,0x10               sub rsp,0x10
sub rsp,0x10{  EC
               10
               89
mov [rbp-0x4],edi { 7D
               FC
```

# Linear Disassembly

## Issues

- ## With ELF files in x86, linear disassembly tends to be useful
  - Compilers do not emit inline data and the process rapidly synchronizes
  - Still, padding and alignment efforts may create some wrong instructions

- ## With PE files, compilers may emit in inline data and Linear Disassembly is not adequate
  - Every time data is found, disassembly becomes desynchronized

- ## Other architectures (ARM) and binary objects usually are not suited for Linear Disassembly
  - Obfuscation may include code as data, which is loaded dynamically
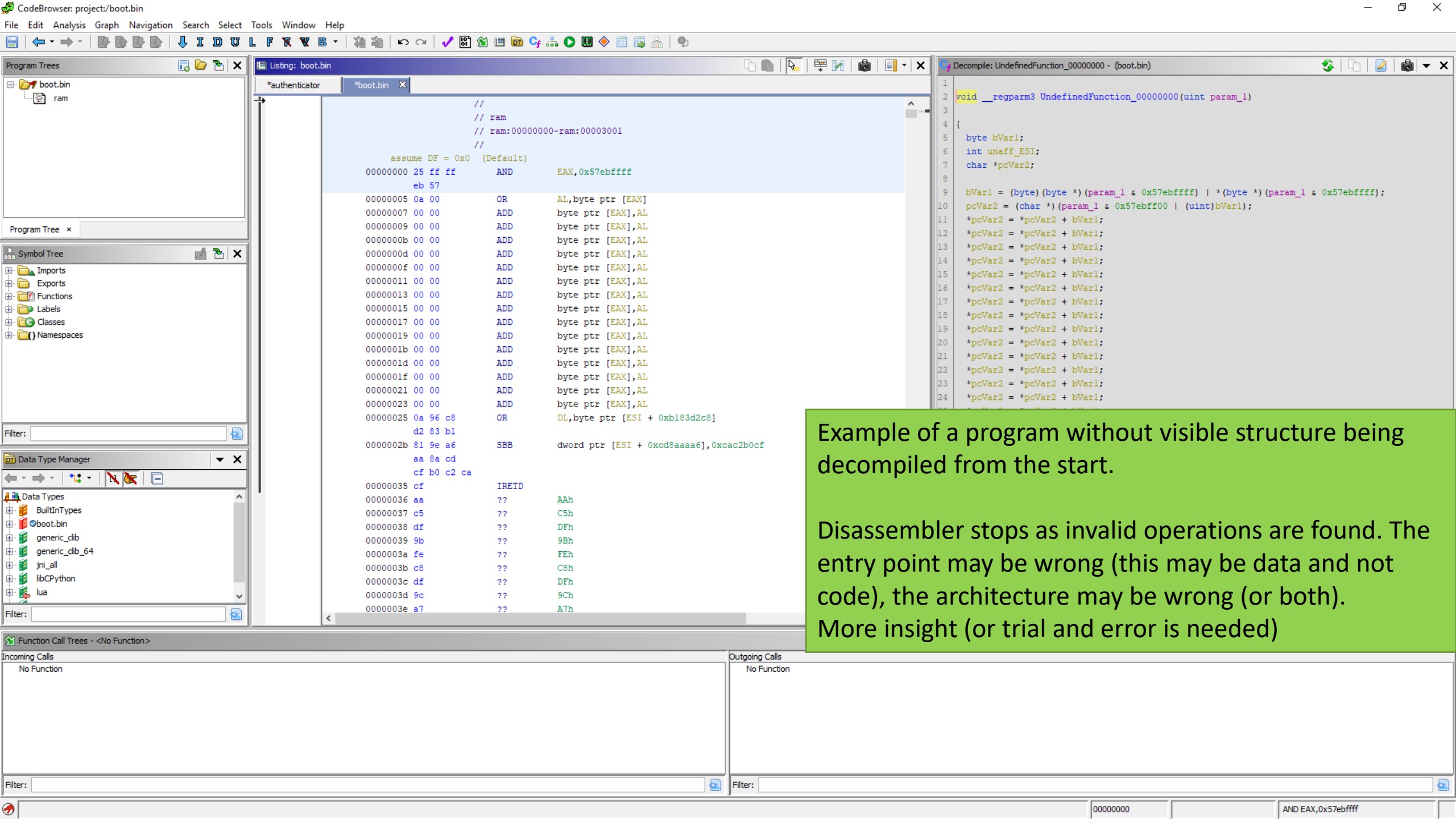  - Fixed length instruction sets will not easily synchronize

# Linear Disassembly

**So why is it useful?**

- Code in the binary blob may be executed with a dynamic call
  - Some JMP/CALL with an address computed dynamically and unknown to the static analyzer

- Linear Disassembly will decompile everything:
  - whether or not it is called - May be useful to uncover hidden program code
  - even if the binary blob is not a structured executable – Boot sector, firmware

- Readily available with simple tools: objdump and gdb
  - Gdb memory dump (x/i) will also use Linear Disassembly

# Recursive Disassembly

- More complex approach that disassembles code since an initial point, **while following the control flow**.
  - That is: follows `jmp`, `call` and `ret`

- **As long as the start point is correct**, or it synchronizes rapidly, flow can be fully recovered
  - This is the standard process for more complex tools such as **ghidra** and **IDA**

- Goes around inline data as no instruction will exist that will make the program execute at such address
  - Well... control flow can easily be forged with `((void (*)(int, char*)) ptr)()`

Example of a program without visible structure being decompiled from the start.

Disassembler stops as invalid operations are found. The entry point may be wrong (this may be data and not code), the architecture may be wrong (or both). More insight (or trial and error is needed)

# Function detection

- Functions frequently include known **prolog** and **epilogues**
  - Prolog: setup the stack and optionally setup Stack Guard Canaries
  - Epilog: optionally check the canaries and release stack

- This information may be used to determine function boundaries
  - But it is architecture and compiler dependent

- Alternatives:
  - Pattern matching (automatic, done by disassemblers) can also recover functions
  - Exception handling code in the .eh_frame section
  - gcc intrinsics to cleanup stacks with exceptions __attribute__((__cleanup__(f))) and __builtin_return_address(n)

# Function detection

## Typical Prologue with Stack Guard

Stack allocation code
- Stores RBP
- Makes RBP = RSP
- Allocates 0x30 bytes

Stores register in stack

```
00400af7 55                 PUSH     RBP
00400af8 48 89 e5           MOV      RBP,RSP
00400afb 48 83 ec 30        SUB      RSP,0x30
00400aff 89 7d dc           MOV      dword ptr [RBP + local_2c],EDI
00400b02 64 48 8b           MOV      RAX,qword ptr FS:[0x28]
         04 25 28
         00 00 00
00400b0b 48 89 45 f8        MOV      qword ptr [RBP + local_10],RAX
00400b0f 31 c0              XOR      EAX,EAX
```

Canary setup
- Fetches value from FS:[0x28] to RAX
- Stores value at RBP+local_10 (top of the local stack)
- Erase RAX

# Function detection

## Typical Epilogue with Stack Guard

Fetches the Canary

```
00400b5a 48 8b 45 f8        MOV        RAX,qword ptr [RBP + local_10]
00400b5e 64 48 33           XOR        RAX,qword ptr FS:[0x28]
         04 25 28
         00 00 00
00400b67 74 05              JZ         LAB_00400b6e
00400b69 e8 b2 fb           CALL       __stack_chk_fail
         ff ff
                    -- Flow Override: CALL_RETURN (CALL_TERMINATOR)

                    LAB_00400b6e
00400b6e c9                 LEAVE
00400b6f c3                 RET
```

- XORs the Canary with reference value
- This sets the Zero flag if they are equal (No corruption)
- Jumps to end of program, or crashes the program with __stack_chk_fail

Deallocate stack and return to caller

# Calling Conventions

- Compilers handle the function calling processes differently, and we have several conventions
  - Adapted to how programmers use the languages (number of arguments)
  - Adapted to number of registers and other architecture details

- These dictate:
  - How arguments are passed to the callee
  - How return codes are passed to the caller
  - Who allocates the stack
  - Who stores important registers such as the Program Counter
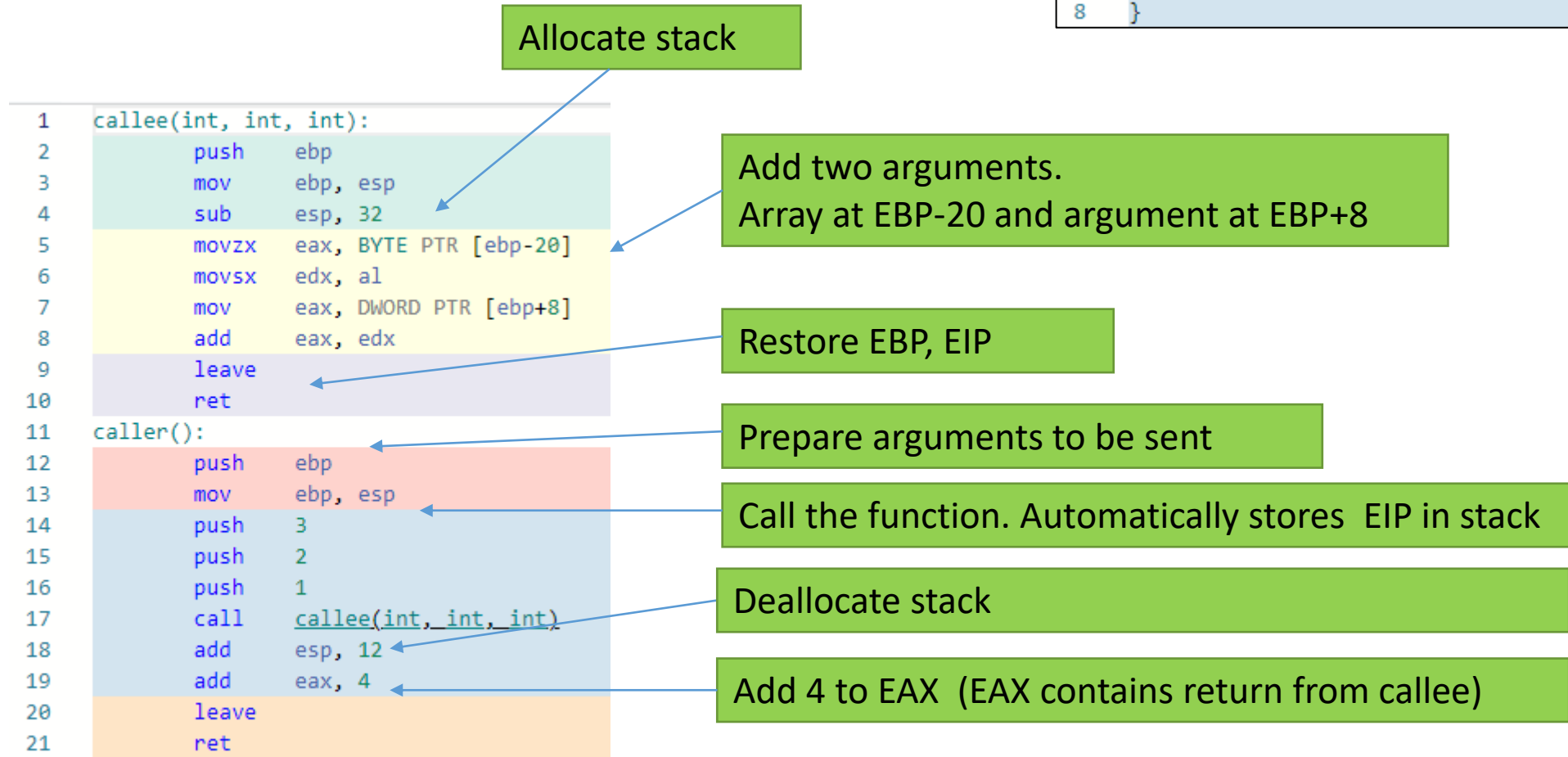
# Calling Conventions

## cdecl

- Originally created by Microsoft compilers, widely used in x86, including GCC
  - Standard method for most code in x86 environments

- **Arguments: passed in the stack, in inverted order (right to left)**
  - First argument is pushed last

- Registers: Mixed
  - Caller saves RIP, A, C, D
  - Callee saves BP, and others and restores RIP

# Calling Conventions

## cdecl

```
1  int callee(int a, int b, int c) {
2      char d[20];
3      return a + d[0];
4  }
5
6  int caller(void) {
7      return callee(1, 2, 3) + 4;
8  }
```

Allocate stack

```
1   callee(int, int, int):
2           push    ebp
3           mov     ebp, esp
4           sub     esp, 32
5           movzx   eax, BYTE PTR [ebp-20]
6           movsx   edx, al
7           mov     eax, DWORD PTR [ebp+8]
8           add     eax, edx
9           leave
10          ret
11  caller():
12          push    ebp
13          mov     ebp, esp
14          push    3
15          push    2
16          push    1
17          call    callee(int,_int,_int)
18          add     esp, 12
19          add     eax, 4
20          leave
21          ret
```

Add two arguments.
Array at EBP-20 and argument at EBP+8

Restore EBP, EIP

Prepare arguments to be sent

Call the function. Automatically stores  EIP in stack

Deallocate stack

Add 4 to EAX  (EAX contains return from callee)
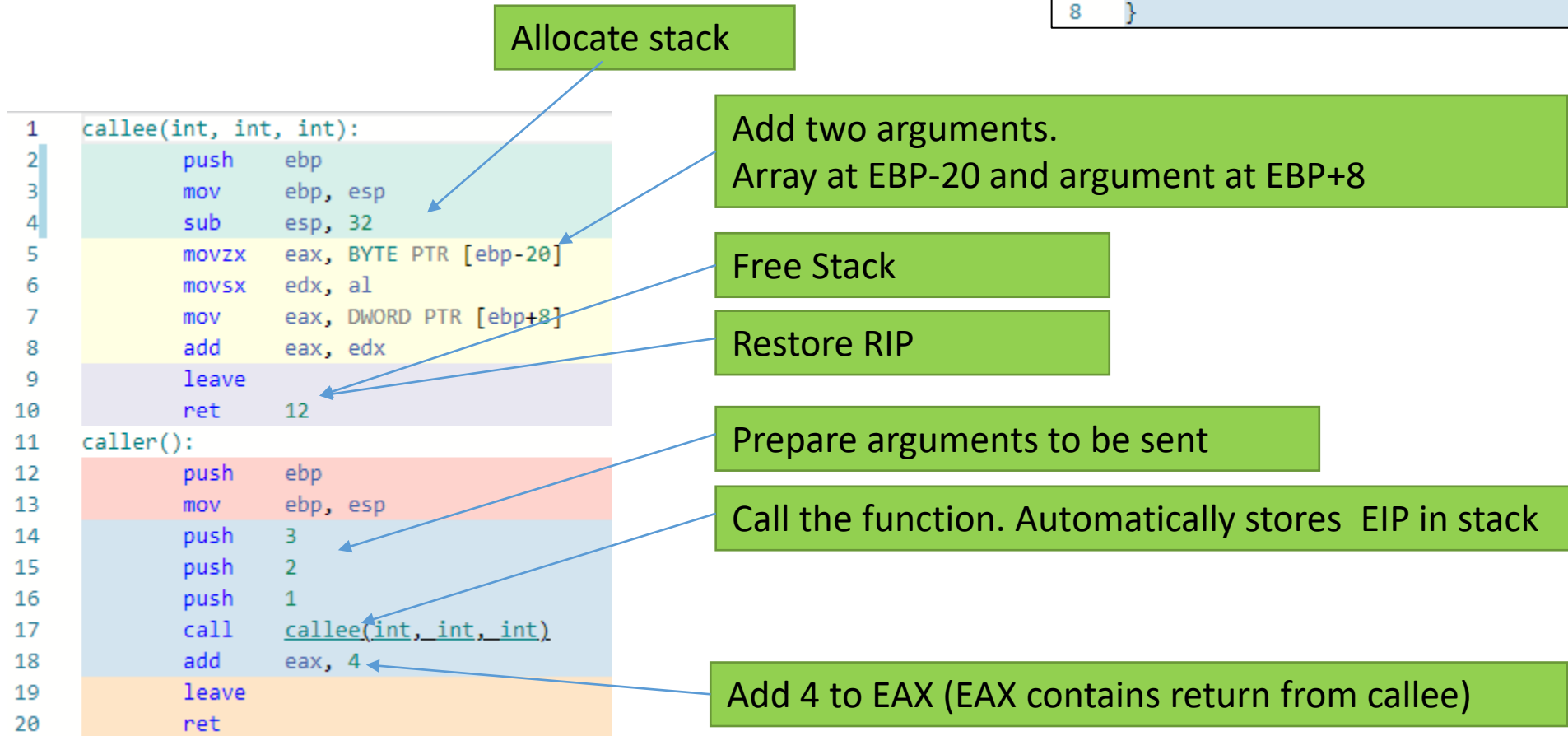
# Calling Conventions

## stdcall

- Official call convention for the Win32API (32 bits)

- **Arguments: passed in the stack from right to left**
  - Additional arguments are passed in the stack

- Registers: Callee saves
  - Except EAX, ECX and EDX which can be freely used

- Stack Red Zone: Leaf functions have a 128 byte area kept safe which doesn't need to be allocated
  - Can be used for local variables, and avoids the use of two operations (sub rsp, add rsp)
  - Leaf functions are functions that do not call others

# Calling Conventions

## stdcall

```
1   int callee(int a, int b, int c) {
2       char d[20];
3       return a + d[0];
4   }
5
6   int caller(void) {
7       return callee(1, 2, 3) + 4;
8   }
```
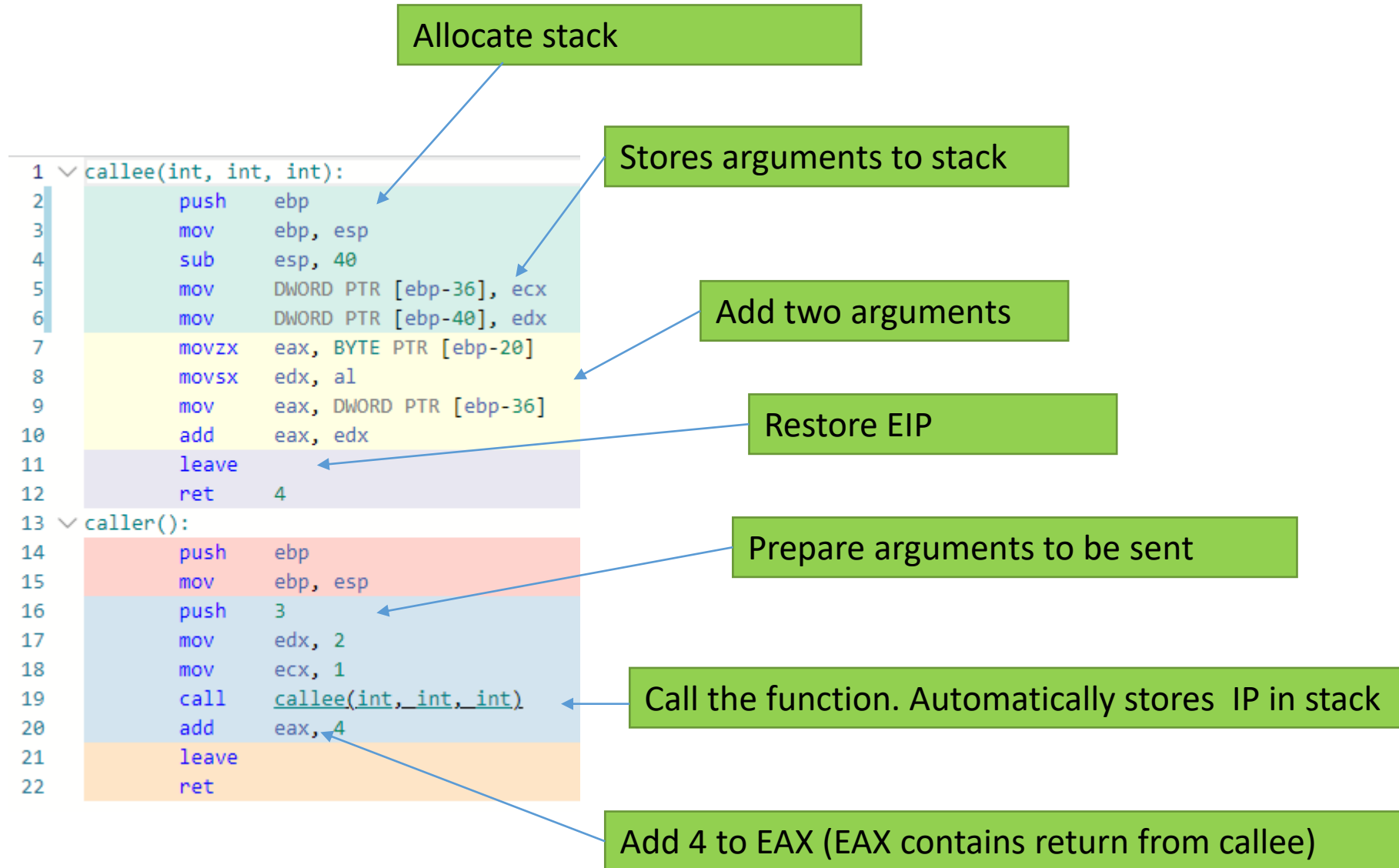
**Allocate stack**

```
1   callee(int, int, int):
2           push    ebp
3           mov     ebp, esp
4           sub     esp, 32
5           movzx   eax, BYTE PTR [ebp-20]
6           movsx   edx, al
7           mov     eax, DWORD PTR [ebp+8]
8           add     eax, edx
9           leave
10          ret     12
11  caller():
12          push    ebp
13          mov     ebp, esp
14          push    3
15          push    2
16          push    1
17          call    callee(int,_int,_int)
18          add     eax, 4
19          leave
20          ret
```

**Add two arguments.**
**Array at EBP-20 and argument at EBP+8**

**Free Stack**

**Restore RIP**

**Prepare arguments to be sent**

**Call the function. Automatically stores  EIP in stack**

**Add 4 to EAX (EAX contains return from callee)**

# Calling Conventions

## fastcall

- Official call convention for Win32API  64bits

- **Arguments: left to right, first as registers**
    - Additional arguments are passed in the stack

- Registers: Caller saves

- Stack Shadow Zone: Leaf functions have a 32 byte area kept safe which doesn't need to be allocated
    - Can be used for local variables, and avoids the use of two operations (sub rsp, add rsp)
    - Leaf functions are functions that do not call others

# Calling Conventions

## fastcall (32bits)

Allocate stack

Stores arguments to stack

Add two arguments

Restore EIP

```
1 ∨ callee(int, int, int):
2        push     ebp
3        mov      ebp, esp
4        sub      esp, 40
5        mov      DWORD PTR [ebp-36], ecx
6        mov      DWORD PTR [ebp-40], edx
7        movzx    eax, BYTE PTR [ebp-20]
8        movsx    edx, al
9        mov      eax, DWORD PTR [ebp-36]
10       add      eax, edx
11       leave
12       ret      4
13 ∨ caller():
14       push     ebp
15       mov      ebp, esp
16       push     3
17       mov      edx, 2
18       mov      ecx, 1
19       call     callee(int,_int,_int)
20       add      eax, 4
21       leave
22       ret
```

Prepare arguments to be sent

Call the function. Automatically stores  IP in stack

Add 4 to EAX (EAX contains return from callee)

# Calling Conventions

```
 1  callee(int, int, int):
 2          push    ebp
 3          mov     ebp, esp
 4          sub     esp, 32
 5          movzx   eax, BYTE PTR [ebp-20]
 6          movsx   edx, al
 7          mov     eax, DWORD PTR [ebp+8]
 8          add     eax, edx
 9          leave
10          ret
11  caller():
12          push    ebp
13          mov     ebp, esp
14          push    3
15          push    2
16          push    1
17          call    callee(int, int, int)
18          add     esp, 12
19          add     eax, 4
20          leave
21          ret
```

**cdecl**

```
 1  callee(int, int, int):
 2          push    ebp
 3          mov     ebp, esp
 4          sub     esp, 32
 5          movzx   eax, BYTE PTR [ebp-20]
 6          movsx   edx, al
 7          mov     eax, DWORD PTR [ebp+8]
 8          add     eax, edx
 9          leave
10          ret     12
11  caller():
12          push    ebp
13          mov     ebp, esp
14          push    3
15          push    2
16          push    1
17          call    callee(int, int, int)
18          add     eax, 4
19          leave
20          ret
```

**stdcall**

```
 1  callee(int, int, int):
 2          push    ebp
 3          mov     ebp, esp
 4          sub     esp, 40
 5          mov     DWORD PTR [ebp-36], ecx
 6          mov     DWORD PTR [ebp-40], edx
 7          movzx   eax, BYTE PTR [ebp-20]
 8          movsx   edx, al
 9          mov     eax, DWORD PTR [ebp-36]
10          add     eax, edx
11          leave
12          ret     4
13  caller():
14          push    ebp
15          mov     ebp, esp
16          push    3
17          mov     edx, 2
18          mov     ecx, 1
19          call    callee(int, int, int)
20          add     eax, 4
21          leave
22          ret
```

**fastcall**

# Calling Conventions

## Fastcall for 64bits (Windows)

- Official convention for x86_64 architectures with MSVC (Windows)
  - Mandatory if compiling for x86_64 in Windows

- Arguments: passed as RDX, RCX, R8, R9
  - Additional arguments are passed in the stack (right to left)

- Registers: Mixed
  - Caller save: RAX, RCX, RDX, R8, R9, R10, R11
  - Callee save: RBX, RBP, RDI, RSI, RSP, R12, R13, R14, and R15

- Stack Red Zone: Leaf functions have a 32 byte area kept safe, allocated by the callee
  - Can be used to store RDX, RCX, R8, R9
  - (Leaf functions are functions that do not call others)

# Calling Conventions

## fastcall (64bits)

```
1   int callee(int a, int b, int c) {
2       char d[20];
3       return a + d[0];
4   }
5
6   int caller(void) {
7       return callee(1, 2, 3) + 4;
8   }
```

```
5    int callee(int,int,int) PROC                    ; callee
6    $LN3:
7            mov      DWORD PTR [rsp+24], r8d
8            mov      DWORD PTR [rsp+16], edx
9            mov      DWORD PTR [rsp+8], ecx
10           sub      rsp, 40                         ; 00000028H
11           mov      eax, 1
12           imul     rax, rax, 0
13           movsx    eax, BYTE PTR d$[rsp+rax]
14           mov      ecx, DWORD PTR a$[rsp]
15           add      ecx, eax
16           mov      eax, ecx
17           add      rsp, 40                         ; 00000
18           ret      0
19   int callee(int,int,int) ENDP                    ; callee
20
21   int caller(void) PROC                            ; caller
22   $LN3:
23           sub      rsp, 40                         ; 0000
24           mov      r8d, 3
25           mov      edx, 2
26           mov      ecx, 1
27           call     int callee(int,int,int)         ; callee
28           add      eax, 4
29           add      rsp, 40                         ; 00000028H
30           ret      0
31   int caller(void) ENDP                            ; caller
```

Stores arguments to shadow

Add two arguments

Free Stack

Restore RIP

Prepare arguments to be sent

Call the function. Automatically stores  IP in stack

Add 4 to EAX (EAX contains return from callee)

# Calling Conventions

## System V AMD64 ABI

- Official convention for x64 architectures using Linux, BSD, Unix, Windows

- Arguments: passed as RDI, RSI, RDX, RCX, R8, R9
  - Additional arguments are passed in the stack

- Registers: Caller saves
  - Except RBX, RSP, RBP, R12-R15 which callee must save if they are used

- Stack Red Zone: Leaf functions have a 128 byte area kept safe which doesn't need to be allocated
  - Can be used for local variables, and avoids the use of two operations (sub rsp, add rsp)
  - Leaf functions are functions that do not call others

# Calling Conventions

## System V AMD64 ABI

```
1  ∨ int callee(int a, int b, int c) {
2        char d[20];
3        return a + d[0];
4    }
5
6  ∨ int caller(void) {
7        return callee(1, 2, 3) + 4;
8    }
```

Leaf function uses stack directly

```
1    callee(int, int, int):
2            push    rbp
3            mov     rbp, rsp
4            mov     DWORD PTR [rbp-36], edi
5            mov     DWORD PTR [rbp-40], esi
6            mov     DWORD PTR [rbp-44], edx
7            movzx   eax, BYTE PTR [rbp-32]
8            movsx   edx, al
9            mov     eax, DWORD PTR [rbp-36]
10           add     eax, edx
11           pop     rbp
12           ret
13   caller():
14           push    rbp
15           mov     rbp, rsp
16           mov     edx, 3
17           mov     esi, 2
18           mov     edi, 1
19           call    callee(int,_int,_int)
20           add     eax, 4
21           pop     rbp
22           ret
```

Adds value with argument

Restore RIP

Prepare arguments to be sent

Call the function. Automatically stores RIP in stack

Add values. EAX will contain the result

# Calling Conventions

## 64bits

```
 1  callee(int, int, int):
 2          push    rbp
 3          mov     rbp, rsp
 4          mov     DWORD PTR [rbp-36], edi
 5          mov     DWORD PTR [rbp-40], esi
 6          mov     DWORD PTR [rbp-44], edx
 7          movzx   eax, BYTE PTR [rbp-32]
 8          movsx   edx, al
 9          mov     eax, DWORD PTR [rbp-36]
10          add     eax, edx
11          pop     rbp
12          ret
13  caller():
14          push    rbp
15          mov     rbp, rsp
16          mov     edx, 3
17          mov     esi, 2
18          mov     edi, 1
19          call    callee(int, int, int)
20          add     eax, 4
21          pop     rbp
22          ret
```

**System V AMD64 ABI**

```
 5  int callee(int,int,int) PROC
 6  $LN3:
 7          mov     DWORD PTR [rsp+24], r8d
 8          mov     DWORD PTR [rsp+16], edx
 9          mov     DWORD PTR [rsp+8], ecx
10          sub     rsp, 40
11          mov     eax, 1
12          imul    rax, rax, 0
13          movsx   eax, BYTE PTR d$[rsp+rax]
14          mov     ecx, DWORD PTR a$[rsp]
15          add     ecx, eax
16          mov     eax, ecx
17          add     rsp, 40
18          ret     0
19  int callee(int,int,int) ENDP
20
21  int caller(void) PROC
22  $LN3:
23          sub     rsp, 40
24          mov     r8d, 3
25          mov     edx, 2
26          mov     ecx, 1
27          call    int callee(int,int,int)
28          add     eax, 4
29          add     rsp, 40
30          ret     0
31  int caller(void) ENDP
```

**fastcall**

# Common Logic Structures

- When analyzing code, it's important to recognize basic flow control structures
  - Remember that the decompiler may be unreliable

- Basic structures:
  - If else
  - Switch case
  - For

# Common Logic Structures

**Conditional Branches (if else)**

- Basic control-flow instructions: move execution to a defined address if a condition is true

  ▪ Usually, one condition tested at a time. Complex If/else must be broken

- Assembly code is structed **as a graph with tests and execution** statements (the conditions body)

- x86 and most architectures have inherent support for many types of comparisons.

  ▪ In x86 this is the jXX family of instructions.

# x86 Opcode Structure and Instruction Overview

Left table (first opcode byte rows 0–F, second nibble 0–F):

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | ADD | | | | | | ES PUSH | ES POP | OR | | | | | | CS PUSH | TWO BYTE |
| 1 | ADC | | | | | | SS PUSH | SS POP | SBB | | | | | | DS PUSH | DS POP |
| 2 | AND | | | | | | ES SEGMENT OVERRIDE | DAA | SUB | | | | | | CS SEGMENT OVERRIDE | DAS |
| 3 | XOR | | | | | | SS | AAA | CMP | | | | | | DS | AAS |
| 4 | INC | | | | | | | | DEC | | | | | | | |
| 5 | PUSH | | | | | | | | POP | | | | | | | |
| 6 | PUSHAD | POPAD | BOUND | ARPL | FS | GS | OPERAND SIZE | ADDRESS SIZE | PUSH | IMUL | PUSH | IMUL | INS | | OUTS | |
| 7 | JO | JNO | JB | JNB | JE | JNE | JBE | JA | JS | JNS | JPE | JPO | JL | JGE | JLE | JG |
| 8 | ADD/ADC/AND/XOR OR/SBB/SUB/CMP | | TEST | | XCHG | | MOV REG | | | | MOV SREG | LEA | MOV SREG | POP | | |
| 9 | NOP | XCHG EAX | | | | | | | CWD | CDQ | CALLF | WAIT | PUSHFD | POPFD | SAHF | LAHF |
| A | MOV EAX | | | | MOVS | | CMPS | | TEST | | STOS | | LODS | | SCAS | |
| B | MOV | | | | | | | | | | | | | | | |
| C | SHIFT IMM | | RETN | | LES | LDS | MOV IMM | | ENTER | LEAVE | RETF | | INT3 | INT IMM | INTO | IRETD |
| D | SHIFT 1 ROL/ROR/RCL/RCR/SHL/SHR/SAL/SAR | | SHIFT CL | | AAM | AAD | SALC | XLAT | FPU | | | | | | | |
| E | LOOPNZ | LOOPZ | LOOP | JECXZ | IN IMM | | OUT IMM | | CALL | JMP | JMPF | JMP SHORT | IN DX | | OUT DX | |
| F | LOCK EXCLUSIVE ACCESS | ICE BP | REPNE | REPE CONDITIONAL REPETITION | HLT | CMC | TEST/NOT/NEG [i]MUL/[i]DIV | | CLC | STC | CLI | STI | CLD | STD | INC DEC | INC/DEC CALL/JMP PUSH |

Right table (two-byte 0F opcodes):

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | {L,S}LDT {L,S}TR VER{R,W} | {L,S}GDT {L,S}IDT {L,S}MSW | LAR | LSL | | CLTS | | INVD | WBINVD | | UD2 | | | NOP | | |
| 1 | SSE{1,2,3} | | | | | | Prefetch SSE1 | | HINT_NOP | | | | | | | |
| 2 | MOV CR/DR | | | | | | SSE{1,2} | | | | | | | | | |
| 3 | WRMSR | RDTSC | RDMSR | RDPMC | SYSENTER | SYSEXIT | | GETSEC SMX | MOVBE / THREE BYTE | THREE BYTE SSE4 | | | | | | |
| 4 | CMOV | | | | | | | | | | | | | | | |
| 5 | SSE{1,2} | | | | | | | | | | | | | | | |
| 6 | MMX, SSE2 | | | | | | | | | | | | | | | |
| 7 | MMX, SSE{1,2,3}, VMX | | | | | | | | | | MMX, SSE{2,3} | | | | | |
| 8 | JO | JNO | JB | JNB | JE | JNE | JBE | JA | JS | JNS | JPE | JPO | JL | JGE | JLE | JG |
| 9 | SETO | SETNO | SETB | SETNB | SETE | SETNE | SETBE | SETA | SETS | SETNS | SETPE | SETPO | SETL | SETGE | SETLE | SETG |
| A | PUSH FS | POP FS | CPUID | BT | SHLD | | | | PUSH GS | POP GS | RSM | BTS | SHRD | | *FENCE | IMUL |
| B | CMPXCHG | | LSS | BTR | LFS | LGS | MOVZX | | POPCNT | UD | BT BTS BTR BTC | BTC | BSF | BSR | MOVSX | |
| C | XADD | | SSE{1,2} | | | | CMPXCHG | | BSWAP | | | | | | | |
| D | MMX, SSE{1,2,3} | | | | | | | | | | | | | | | |
| E | MMX, SSE{1,2} | | | | | | | | | | | | | | | |
| F | MMX, SSE{1,2,3} | | | | | | | | | | | | | | | |

## Legend

- **Arithmetic & Logic** (green)
- **Memory** (red)
- **Stack** (orange)
- **Control Flow & Conditional** (blue)
- **Prefix** (yellow)
- **System & I/O** (purple)
- **No Operation (NOP) / Multiple Instructions / Extended Instruction Set** (grey)

## General Opcode Structure

| Element Information # of bytes | Prefix 0-4 | Opcode 1-3 | AddrMode (mod, reg, r/m) 0-1 | SIB Byte (scale, index, base) 0-1 | Displacement 0/1/2/4 | Immediate Data 0/1/2/4 |
|---|---|---|---|---|---|---|

Bit structure:
- Main Opcode bits
- Direction bit
- Operand length bit
- r/m field
- Register/Opcode modifier, defined by primary opcode
- Addressing mode
- Base field
- Index field
- Scale field

## Addressing Modes

| mod | 00 | | 01 | | 10 | | 11 |
|---|---|---|---|---|---|---|---|
| r/m | 16bit | 32bit | 16bit | 32bit | 16bit | 32bit | r/m // REG |
| 000 | [BX+SI] | [EAX] | [BX+SI]+disp8 | [EAX]+disp8 | [BX+SI]+disp16 | [EAX]+disp32 | AL / AX / EAX |
| 001 | [BX+DI] | [ECX] | [BX+DI]+disp8 | [ECX]+disp8 | [BX+DI]+disp16 | [ECX]+disp32 | CL / CX / ECX |
| 010 | [BP+SI] | [EDX] | [BP+SI]+disp8 | [EDX]+disp8 | [BP+SI]+disp16 | [EDX]+disp32 | DL / DX / EDX |
| 011 | [BP+DI] | [EBX] | [BP+DI]+disp8 | [EBX]+disp8 | [BP+DI]+disp16 | [EBX]+disp32 | BL / BX / EBX |
| 100 | [SI] | SIB | [SI]+disp8 | SIB+disp8 | [SI]+disp16 | SIB+disp32 | AH / SP / ESP |
| 101 | [DI] | disp32 | [DI]+disp8 | [EBP]+disp8 | [DI]+disp16 | [EBP]+disp32 | CH / BP / EBP |
| 110 | disp16 | [ESI] | [BP]+disp8 | [ESI]+disp8 | [BP]+disp16 | [ESI]+disp32 | DH / SI / ESI |
| 111 | [BX] | [EDI] | [BX]+disp8 | [EDI]+disp8 | [BX]+disp16 | [EDI]+disp32 | BH / DI / EDI |

## SIB Byte Structure

| encoding | scale (2bit) | Index (3bit) | Base (3bit) |
|---|---|---|---|
| 000 | $2^0$=1 | [EAX] | EAX |
| 001 | $2^1$=2 | [ECX] | ECX |
| 010 | $2^2$=4 | [EDX] | EDX |
| 011 | $2^3$=8 | [EBX] | EBX |
| 100 | – | none | ESP |
| 101 | – | [EBP] | disp32 + disp8 + [EBP] / disp32 + [EBP] |
| 110 | – | [ESI] | ESI |
| 111 | – | [EDI] | EDI |

SIB value = index * scale + base

# Common Logic Structures

**Conditional Branches (if else)**

- Signed comparison: **l < , le <=, g >, ge >=**

- Unsigned comparison: **b <, be <=, a >=, ae >=**
  - Below and Above

- Equality **e**

- Every condition can be negated with **n**

# Common Logic Structures

## Conditional Branches (if else)

- **z, s, c, o,** and **p** for **ZF, SF, CF OF,** and **PF**
  - **ZF: Zero Flag**, 1 if last operation was 0
  - **CF: Carry Flag**. Last operation required an additional bit (e.g. 255 + 1, which has 9 bits)
  - **OF: Overflow Flag**. Last operation had an arithmetic overflow (127 + 127 in a signed variable results in overflow)
  - **PF: Parity Flag**. 1 if last operation resulted in a value with even number of 1
  - **SF: Sign Flag**. 1 if last operation resulted in a signed value (MSB bit = 1)

- **s** means negative, **ns** non-negative
  - Signal or not signal

- **p** and **np** are also pe "parity even" and **po** "parity odd"

# Common Logic Structures

## Conditional Branches (if else)

- **and, or,** and **xor** clear **OF** and **CF**, and set **ZF**, **SF**, and **PF** based on the result

- **test** is like **and** but only sets the flags discarding the result

- Checking **nz** after **test** is like **if (x & mask) in C**

- **test** a register against itself is the fastest way to check for zero or negative

# Common Logic Structures

## Conditional Branches (if else)

- **Direct jump**: target(s) specified in code (harcoded)

- **Indirect jump**: target selected from runtime data like register or memory contents

- **Conditional jump**: target differs based on a condition

# Common Logic Structures

## Conditional Branches (If else)

- Structure can be recognized by one or more conditional branches, without loops

- je: jump equal

- js: jump is sign

- …etc…

```
1  int bar(int b) {
2      return b * b;
3  }
4
5  int foo(int a) {
6      if(a == 0){
7          return bar(a) * 1;
8
9      }
10     else
11         if(a < 0){
12             return bar(a) - 1;
13         }
14         else{
15             return bar(a) + 1;
16         }
17 }
```

```
1  ∨ bar:
2          imul    edi, edi
3          mov     eax, edi
4          ret
5  ∨ foo:
6          test    edi, edi
7          je      .L6
8          js      .L7
9          call    bar
10         add     eax, 1
11         ret
12 ∨ .L6:
13         call    bar
14         ret
15 ∨ .L7:
16         call    bar
17         sub     eax, 1
18         ret
```

# Common Logic Structures

## Switch case

- Structure can be recognized by several comparisons and jumps or **jump table**

- Observe the difference between what a programmer writes and what is produced
  - Switch is written as an atomic instruction, but it isn't
  - Also, it is dangerous because of missing breaks;

- Test: compare two registers. Set 3 flags:
  - PF: Even number of bits
  - ZF: Zero
  - SF: Signed value

```
1   int bar(int b) {
2       return b * b;
3   }
4
5   int foo(int a) {
6       switch(a){
7           case 0:
8               a = bar(1) + 1;
9               break;
10          case 1:
11              a = bar(2+ a) + 2;
12              break;
13          case 3:
14              a = bar(3) + 3;
15          default:
16              a= bar(4) + 4;
17      }
18
19      return a;
20  }
```

```
1   bar:
2           imul    edi, edi
3           mov     eax, edi
4           ret
5   foo:
6           test    edi, edi
7           je      .L3
8           cmp     edi, 1
9           je      .L4
10          mov     edi, 4
11          call    bar
12          add     eax, 4
13          ret
14  .L3:
15          mov     edi, 1
16          call    bar
17          add     eax, 1
18          ret
19  .L4:
20          add     edi, 2
21          call    bar
22          add     eax, 2
23          ret
```

# Common Logic Structures

## loops

- For, while and do while are generally the same

- Identified by:
  - an index
  - an increment
  - a comparison
  - two jumps

```
1    int bar(int b) {
2        return b * b;
3    }
4
5    int foo(int a) {
6        int b = 0;
7        for(int i = 0; i< a; i++){
8            b += bar(i);
9        }
10
11       return b;
12   }
13
14   int caller(void) {
15       return callee(1, 2, 3) + 4;
16   }
```

```
1   bar:
2           imul    edi, edi
3           mov     eax, edi
4           ret
5   foo:
6           push    r12
7           push    rbp
8           push    rbx
9           mov     r12d, edi
10          mov     ebx, 0
11          mov     ebp, 0
12  .L3:
13          cmp     ebx, r12d
14          jge     .L6
15          mov     edi, ebx
16          call    bar
17          add     ebp, eax
18          add     ebx, 1
19          jmp     .L3
20  .L6:
21          mov     eax, ebp
22          pop     rbx
23          pop     rbp
24          pop     r12
25          ret
```

Prepares stack

- r12d will contain the number of iterations
- ebx will be the counter

- Loop body

Jump to top of loop

# C++ code

- C++ is very popular, and adds an additional layer of complexity
  - A program doesn't have functions, has methods

  - Methods have a shared context (the object)

  - Methods can be overridden due to inheritance

  - The `this` pointer commonly allows access to data outside the function stack

  - Contructors, new…?

  - Strings are complex objects

# C++ code

- this pointer
  - The "this" pointer plays a crucial role in the identification of C++ sections in the assembly code. It is initialized to point to the object used, to invoke the function, when it is available in non-static C++ functions.

- Vtables
  - Eases runtime resolution of calls to virtual functions.
  - The compiler generates a vtable containing pointers to each virtual function for the classes which contain virtual functions.

- Constructors and destructors
  - A member function which initializes objects of a class and it can be identified in assembly by studying the objects in which it's created.

# C++ code

- ## Runtime Type Information (RTTI)
  - Mechanism to identify the object type at run.
  - These keywords pass information, such as class name and hierarchy, to the class.

- ## Structured exception handling (SEH)
  - Irregularities in source code that unexpectedly strike during runtime, terminating the program.
  - SEH is the mechanism that controls the flow of execution and handles errors by isolating the code section where the unexpected condition originates. Inheritance

- ## Inheritance
  - allows new objects to take on existing object properties.
  - Observing RTTI relationships can reveal inheritance hierarchy

# hello1.cpp

## A simple hello world

```cpp
1   #include <iostream>
2   #include <string>
3
4   class A {
5       std::string text1;
6
7       public:
8       A(std::string text1) {
9           this->text1 = text1;
10      }
11      void print() {
12          std::cout << this->text1 << std::endl;
13      }
14  };
15
16  int main(int argc, char** argv) {
17      A a(std::string("Hello World"));
18      a.print();
19
20  }
```

# hello1.cpp

```
1   $ readelf --dyn-sym hello1
2
3   Symbol table '.dynsym' contains 21 entries:
4      Num:    Value          Size Type    Bind    Vis      Ndx Name
5        0: 0000000000000000     0 NOTYPE  LOCAL   DEFAULT  UND
6        1: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSt7__cxx1112basic_stri@GLIBCXX_3.4.21 (2)
7        2: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZSt4endlIcSt11char_trait@GLIBCXX_3.4 (4)
8        3: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSt7__cxx1112basic_stri@GLIBCXX_3.4.21 (2)
9        4: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND __cxa_atexit@GLIBC_2.2.5 (3)
10       5: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZStlsIcSt11char_traitsIc@GLIBCXX_3.4.21 (2)
11       6: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSolsEPFRSoS_E@GLIBCXX_3.4 (4)
12       7: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSaIcED1Ev@GLIBCXX_3.4 (4)
13       8: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSt7__cxx1112basic_stri@GLIBCXX_3.4.21 (2)
14       9: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSt7__cxx1112basic_stri@GLIBCXX_3.4.21 (2)
15      10: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSt8ios_base4InitC1Ev@GLIBCXX_3.4 (4)
16      11: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND __gxx_personality_v0@CXXABI_1.3 (5)
17      12: 0000000000000000     0 NOTYPE  WEAK    DEFAULT  UND _ITM_deregisterTMCloneTab
18      13: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _Unwind_Resume@GCC_3.0 (6)
19      14: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSaIcEC1Ev@GLIBCXX_3.4 (4)
20      15: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND __libc_start_main@GLIBC_2.2.5 (3)
21      16: 0000000000000000     0 NOTYPE  WEAK    DEFAULT  UND __gmon_start__
22      17: 0000000000000000     0 NOTYPE  WEAK    DEFAULT  UND _ITM_registerTMCloneTable
23      18: 0000000000000000     0 FUNC    GLOBAL  DEFAULT  UND _ZNSt8ios_base4InitD1Ev@GLIBCXX_3.4 (4)
24      19: 0000000000000000     0 FUNC    WEAK    DEFAULT  UND __cxa_finalize@GLIBC_2.2.5 (3)
25      20: 00000000000040a0   272 OBJECT  GLOBAL  DEFAULT   26 _ZSt4cout@GLIBCXX_3.4 (4)
```

# C++ code

No C++ class declarations, but C++ class use.

- Constructors
- Methods
- Destructors

```
1
2   /* WARNING: Unknown calling convention yet parameter storage is locked */
3
4   int main(void)
5
6   {
7     A local_68 [32];
8     basic_string<char,std::char_traits<char>,std::allocator<char>> local_48 [47];
9     allocator<char> local_19 [9];
10
11    std::allocator<char>::allocator();
12                    /* try { // try from 00101203 to 00101207 has its CatchHandler @ 00101263 */
13    std::__cxx11::basic_string<char,std::char_traits<char>,std::allocator<char>>::basic_string
14           ((char *)local_48,(allocator *)"Hello World");
15                    /* try { // try from 00101216 to 0010121a has its CatchHandler @ 00101252 */
16    A::A(local_68,(basic_string)0xb8);
17    std::__cxx11::basic_string<char,std::char_traits<char>,std::allocator<char>>::~basic_string
18           (local_48);
19    std::allocator<char>::~allocator(local_19);
20                    /* try { // try from 0010123a to 0010123e has its CatchHandler @ 0010127d */
21    A::print(local_68);
22    A::~A(local_68);
23    return 0;
24  }
25
```

# C++ code



```
00101216 e8 e1 00        CALL        A::A
         00 00
                         } // end try from 00101216 to 0010121a
0010121b 48 8d 45 c0      LEA         RAX=>local_48,[RBP + -0x40]
0010121f 48 89 c7         MOV         RDI,RAX
00101222 e8 19 fe         CALL        ~basic_string
         ff ff
00101227 48 8d 45 ef      LEA         RAX=>local_19,[RBP + -0x11]
0010122b 48 89 c7         MOV         RDI,RAX
0010122e e8 4d fe         CALL        ~allocator
         ff ff
00101233 48 8d 45 a0      LEA         RAX=>local_68,[RBP + -0x60]
00101237 48 89 c7         MOV         RDI,RAX

                         try { // try from 0010123a to 0010123e has its CatchHandler @...
         LAB_0010123a                              XREF[1]:    00]
0010123a e8 11 01         CALL        A::print
         00 00
                         } // end try from 0010123a to 0010123e
0010123f 48 8d 45 a0      LEA         RAX=>local_68,[RBP + -0x60]
00101243 48 89 c7         MOV         RDI,RAX
00101246 e8 3d 01         CALL        A::~A
         00 00
0010124b b8 00 00         MOV         EAX,0x0
         00 00
00101250 eb 45            JMP         LAB_00101297
```

Additional Hints related to exception handling

Standard ASM code with function invocation, using arguments in registers and values stored in the stack

# C++ code

.eh_frame ELF section contains information about the multiple methods.

Required for unwinding frames, when iterating over the function frames. Contains language specific information, organized in Call Frame Information records

```
***************************************************************************
* Frame Descriptor Entry                                                  *
***************************************************************************
                  fde_00102148                          XREF[1]:    0010205c(*)
00102148 1c 00 00 00     ddw        1Ch                          (FDE) Length
0010214c a4 00 00 00     ddw        cie_001020a8                 (FDE) CIE Reference Pointer
00102150 00 f2 ff ff     ddw        A::print                     (FDE) PcBegin
00102154 37 00 00 00     ddw        37h                          (FDE) PcRange
00102158 00              uleb128    0h                           (FDE) Augmentation Data Length
00102159 41 0e 10        db[15]                                  (FDE) Call Frame Instructions
         86 02 43
         0d 06 72 ...
```

# C++ code

this is passed as an additional, hidden argument
In this case, in RDI as the method has no arguments

```
*************************************************************
* A::print()                                                *
*************************************************************
                 undefined __thiscall print(A * this)
    undefined         AL:1           <RETURN>
    A *               RDI:8 (auto)   this
    undefined8        Stack[-0x10]:8 local_10                      XREF[2]:     00101358(W),
                                                                                0010135c(R)
                 _ZN1A5printEv                                XREF[4]:    Entry Point(*), main:0010123a(c),
                 A::print                                                 00102058, 00102150(*)
00101350 55                 PUSH       RBP
00101351 48 89 e5           MOV        RBP,RSP
00101354 48 83 ec 10        SUB        RSP,0x10
00101358 48 89 7d f8        MOV        qword ptr [RBP + local_10],this
0010135c 48 8b 45 f8        MOV        RAX,qword ptr [RBP + local_10]
00101360 48 89 c6           MOV        RSI,RAX
00101363 48 8d 3d           LEA        this,[std::cout]                        =
         36 2d 00 00
0010136a e8 f1 fc           CALL       operator<<                      basic_ostream * operator<<(basic...
         ff ff
0010136f 48 89 c2           MOV        RDX,RAX
00101372 48 8b 05           MOV        RAX,qword ptr [->endl<char,std::char_traits<ch... = 00105008
         57 2c 00 00
00101379 48 89 c6           MOV        RSI=>endl<char,std::char_traits<char>>,RAX      = ??
0010137c 48 89 d7           MOV        this,RDX
0010137f e8 ec fc           CALL       operator<<                      undefined operator<<(basic_ostre...
         ff ff
```