Automating Binary Analysis with Ghidra's P-Code

Gergely Revay

SIEMENS

Unrestricted | © Siemens 2021 | Gergely Revay | T RDA CST SES | 2021-10-07

Gergely Revay Senior Security Researcher

LinkedIn: <u>https://www.linkedin.com/in/gergelyrevay/</u> Twitter: @geri_revay

Youtube: <u>https://youtube.com/aetherlabnet</u>



Unrestricted | © Siemens 2021 | Gergely Revay | T RDA CST SES | 2021-10-07

Agenda

- Ghidra
- Ghidra Scripts
 - Intro
 - Python vs Java
 - Headless Mode
 - Flat API vs SDK
- P-Code
 - Intro
 - Raw P-Code
 - High P-Code
 - But Why?
 - Examples
- Recap
- QʻnʻA



https://media.makeameme.org/created/agenda-5b5740.jpg



What is Ghidra

NSA releases Ghidra, a free software reverse engineering toolkit

NSA's Ghidra greeted with positive reviews by the infosec community.



https://www.zdnet.com/article/nsa-release-ghidra-a-free-software-reverse-engineering-toolkit/



	← → → -	🛛 🔁 💽		IDU	LFK	🕷 🖪 🔹	22	🗸 🛱 🖄 🛅 🛅	😋 🚠 🜔 🛄 🔶 🗐 📑 🚠 🛛 🌒	
--	---------	-------	--	-----	-----	-------	-----------	-----------	---------------------	--

Program Trees	🖥 🔁 🔁	🖽 Listing: zyppe_7 - (68	addresses selected)		Po 💼 🏊 🗮 👫 💩 🗐 🗸	×	Decompile: encrypt - (zyppe_7)	🌮 🔻 🕲 🖓 🖓
🖃 🏹 zyppe_7	^	*zyppe_7 🗙				23	int bul	^
.bss			00401702 01 01 51	00	THD_00401043	24	int local 18:	
.data			00 00 00		,	25	int i:	
			004017b8 8b 45 f4	MOV	EAX, dword ptr [RBP + j]	26	int swapper]:	
jot			004017bb 48 98	CDQE		27	int i:	
			004017bd 8b 94 85	MOV	EDX, dword ptr [RBP + RAX*0x4 + -0x460]	28		
- 🗟 .fini_array			a0 fb ff ff			29	local 68 = 0x7465726365732041;	
			004017c4 8b 45 f8	MOV	EAX, dword ptr [RBP + swapper1]	30	$\log_{10} f_0 = 0 \times 6 c_{20} f_{16} f_{16} c_{20} c_$	
			004017c7 8d 34 02	LEA	ESI,[RDX + RAX*0x1]	31	local 58 = $0x2061207265676e6f$;	
😰 .eh_frame			004017ca 8b 4d f4	MOV	ECX, dword ptr [RBP + j]	32	$\log_{10} \log_{10} $	
eh_frame_hdr			004017cd ba 4f ec	MOV	EDX,0x4ec4ec4f	33	local 48 = 0x656d6f732065636e:	
.rodata			c4 4e			34	$\log_{10} 40 = 0 x776666b2065666f;$	
- 🗟 .fini			004017d2 89 c8	MOV	EAX, ECX	35	local 38 = 0x74692073:	
.text			004017d4 f7 ea	IMUL	EDX	36	$\log_{10} 34 = 0$:	
🗟 ,plt			004017d6 cl fa 04	SAR	EDX, 0x4	37	i = 0:	
En init	¥		004017d9 89 c8	MOV	EAX, ECX	38	while $(i < 0x100)$ {	
Program Tree × DWARF ×			004017db cl f8 lf	SAR	EAX, 0x1f	39	key base[i] = i:	
			004017de 29 c2	SUB	EDX, EAX	40	i = i + 1;	
Symbol Tree	🗾 🖄 🗙		004017e0 89 d0	MOV	EAX, EDX	41		
🖶 📴 c_str	^		004017e2 6b c0 34	IMUL	EAX,EAX,0x34	42	swapper1 = 0;	
🕀 📴 dose			004017e5 29 cl	SUB	ECX, EAX	43	i = 0;	
i compare			004017e7 89 c8	MOV	EAX, ECX	44	while (j < 0x100) {	
teregister_tm_clones			004017e9 48 98	CDQE		45	<pre>iVar2 = key base[j] + swapper1 + (int)*(char *)((long)slocal 68 + (long)(j % 0x34));</pre>	
the contract of the second second			004017eb 0f b6 44	MOVZX	EAX, byte ptr [RBP + RAX*0x1 + -0x60]	46	uVarl = (uint) (iVar2 >> 0x1f) >> 0x18;	
⊞ frame dummy			05 a0			47	<pre>swapper1 = (iVar2 + uVar1 & 0xff) - uVar1;</pre>	
			004017f0 Of be c0	MOVSX	EAX, AL	48	<pre>swapper2 = key_base[j];</pre>	
🕀 🛅 getenv			p04017f3 8d 14 06	LEA	EDX, [RSI + RAX*0x1]	49	<pre>key_base[j] = key_base[swapper1];</pre>	
🕀 📴 Init	_		00401716 89 d0	MOV	EAX, EDX	50	<pre>key_base[swapper1] = swapper2;</pre>	
🖨 🕈 main			00401765 c1 68 19	SAR	EAA, UXII	51	j = j + 1;	
Incal_20			00401715 01 03	ADD	ERA, UXIO	52	}	
Incal_28			00401800 0f b6 d2	MOUZY	EDV, EAA	53	local_18 = 0;	
B local 58			00401803 29 c2	SUB	EDX, DE	54	<pre>swapper1 = 0;</pre>	
⊕ local_78	~		00401805 89 d0	MOV	FAX. FDX	55	keyl = 0;	
<	>		00401807 89 45 f8	MOV	dword ptr [RBP + swapper]].EAX	56	<pre>plain_text = 0;</pre>	
Filter:			0040180a 8b 45 f4	MOV	EAX, dword ptr [RBP + j]	57	<pre>while (plain_text < 0x400) {</pre>	
			0040180d 48 98	CDQE		58	$uVarl = (uint)(local_18 + 1 >> 0x1t) >> 0x18;$	
Data Type Manager	- ×		0040180f 8b 84 85	MOV	EAX, dword ptr [RBP + RAX*0x4 + -0x460]	59	$10cal_{18} = (10cal_{18} + 1 + uvart \epsilon 0xir) - uvart;$	
🌰 - 🔿 - 🔤 - 📐 🔭			a0 fb ff ff			60	uvari = (unc) (xey_base[local_is] + swapperi >> 0xii) >> 0xis;	
			00401816 89 45 e4	MOV	dword ptr [RBP + swapper2],EAX	62	local 28 - key base[local 18]:	
BuiltInTupos			00401819 8b 45 f8	MOV	EAX,dword ptr [RBP + swapper1]	63	key hase[local 18] = key hase[swapper]].	
Carpone 7			0040181c 48 98	CDQE		64	<pre>key_base[seapper1] = local 28:</pre>	
apperic dib 64			0040181e 8b 94 85	MOV	EDX, dword ptr [RBP + RAX*0x4 + -0x460]	65	<pre>uVarl = (uint) (key base[local 18] + key base[swapper]] >> 0x1f) >> 0x18;</pre>	
ini all			a0 fb ff ff			66	<pre>kev2 = kev base[(int)((kev base[local 18] + kev base[swapperl] + uVarl & 0xff) - uVap</pre>	cl)];
windows vs12 32			00401825 8b 45 f4	MOV	EAX, dword ptr [RBP + j]	67	<pre>param 1[plain text] = param 1[plain text] ^ (byte)key1 ^ (byte)key2;</pre>	-,1.
			00401828 48 98	CDQE		68	key1 = key2;	
			0040182a 89 94 85	MOV	dword ptr [RBP + RAX*0x4 + -0x460],EDX	69	<pre>plain_text = plain_text + 1;</pre>	
			a0 fb ff ff			70	}	
			00401831 86 45 18	MOV	EAX, dword ptr [RBP + swapper1]	71	return;	
		↓ <u>↓</u>	00401834 48 98	CDQE	TDV dward star (DDD), successful	/ 72	}	
		<			>	73		×.
		🖳 Console - Scripting						🗎 🌽 🗙
							Activate V	vindows
							Go to Setting	is to activate Windows.
Filter:								
Address pat found in program mem	orus fffffffffbog						004017F2	LEA EDV [DST + DAV#0v1]

Ghidra Scripts

- Like IDAPython
- Extend the functionality
- Full automation is possible
- Java or Python (Jython)
- FlatAPI vs SDK

```
Youyspublic class Main {=-you reversespublic static String reverseString(String str) {=-> StringBuilder reverse = new StringBuilder();=-> for (int idx = hello.length() - 1; idx >= 0; idx--) {=> public static void main(hello.charAt(idx));=-> a> public static void main(String[] args) {=-> String hello = "Hello world!";=-> System.out.println(reverseString(hello));=-
```

The guy she tells you not to worry about

hello = 'Hello world!'¤ print(hello[::-1])¤¬

https://www.reddit.com/r/ProgrammerHumor/comments/66jj7f/java_vs_python/



Script Manager - 240 script	s				O Q 📰	● ≯ × ==	📕 🧐 🔁 📕	×
Scripts 🔨	In Tool	Status	Name	E Description	Кеу	Category	Modified	
			AddCommentToProgramScript.java			Examples	02/12/2020	~
Analysis			AddCommentToProgramScriptPy.py	Adds a comment to a program. DISCLAIMER		Examples->Pyt	02/12/2020	
			AddReferencesInSwitchTable.java	With cursor on switch's "add pc," command		ARM	02/12/2020	
Assembly			AddSingleReferenceInSwitchTable.java	With a user-inputed base address, this scrip		ARM	02/12/2020	
			AppleSingleDoubleScript.java	Given a raw binary Apple Single/Double ima		Binary	02/12/2020	
			ArmThumbFunctionTableScript.java	Makes functions out of a run of selected AR		ARM	02/12/2020	
Conversion			AsciiToBinaryScript.java	Converts an ascii hex file into binary file. W		Conversion	02/12/2020	
+ CustomerSubmissi			AskScript.java	An example of asking for user input. Note th		Examples	02/12/2020	
Data			AskScriptPy.py	An example of asking for user input. Note th		Examples->Pyt	02/12/2020	
- 🛅 Data Types			AssembleBlockScript.java	Assemble hard-coded block of instructions.		Assembly	02/12/2020	
🗄 🫅 Examples			AssembleCheckDevScript.java	Test assembly of the instruction under the c	Ctrl-H	Assembly	02/12/2020	
- EunctionID			AssembleScript.java	Assemble a single instruction, overwriting th		Assembly	02/12/2020	
- Eunctions			AssemblyThrasherDevScript.java	Thoroughly test the assembler by attemptin		Assembly	02/12/2020	
FunctionStartPatte			AutoRenameLabelsScript.java	Renames default labels in a selected region,		Symbol	02/12/2020	
- ELP			AutoRenameSimpleLabels.java	A ghidra script that renames simple function		Symbol	02/12/2020	
- Images			AutoVersionTrackingScript.java	An example of how to create Version Tracki		Examples->Ver	02/12/2020	
Import			BadInstructionCleanup.java	This script cleans up the disassembly for kex		iOS	02/12/2020	
Instructions			BatchRename.java	Recursively finds a folder that matches a str		Project	02/12/2020	
- IOS			BatchSegregate64bit.java	Separates co-mingled n-bit and 64-bit binari		Project	02/12/2020	
			BinaryToAsciiScript.java	Converts a binary file into an ascii hex file.		Conversion	02/12/2020	
			BTreeAnnotationScript.java	Annotates an HFS+ attributes b-Tree file.		iOS	02/12/2020	
Memory			BuildFuncDB.java			CodeAnalysis	02/12/2020	
MultiUser			BuildGhidraJarScript.java	An example of building a single minimal Ghi		Examples	02/12/2020	
PCode			Call AnotherScript.java	Example of a script calling another script.		Examples->De	02/12/2020	
Processor			CallAnotherScriptForAllPrograms.java	Shows how to run a script on all of the progr		Examples	02/12/2020	
- 🛅 Program			CallAnotherScriptForAllProgramsPy.py	Shows how to run a script on all of the progr		Examples->Pyt	02/12/2020	
Project			CallAnotherScriptPy.py	Example of a script calling another script. DI		Examples->Pyt	02/12/2020	
Doforoncon V			ChangeDataSettingsScript.java	Changes the dsiplay settings of the current		Examples	02/12/2020	
< >>			ChooseDataTypeScript java	Example of a script prompting the user for a		Evamples-NDa	02/12/2020	Y

_

Python vs. Java

THAT'S YOUR CONSCIENCE TALKING

UNISIPA

https://pics.me.me/oh-noone-is-judging-you-sweetie-thats-your-conscience-talking-imgflip-com-no-51807818.png

- Whatever your flavor is, nobody is judging
- Ghidra is written in Java
- Python scripting works pretty well through Jython (python 2)
- Python 3: <u>https://github.com/justfoxing/ghidra_bridge</u>
- Java development environment is great
- Ghidra plugin for Eclipse
- Java dev is supported by NSA
- You might need to let go of your principles



Headless Mode

analyzeHeadless /Users/user/ghidra/projects MyProject -import hello.exe -preScript GetInfoScript.java

- Way to implement fully automated scripts
- Calling it could be nicer



https://sadanduseless.b-cdn.net/wp-content/uploads/2021/07/headless-gymnasts1.jpg



- Interface to write simple scripts
- Reliability in long term is the main goal
- ~147 methods available
- GhidraScript class is a subclass of FlatProgramAPI
- Stick to it if possible

NOTE:

NO METHODS SHOULD EVER BE REMOVED FROM THIS CLASS.
 NO METHOD SIGNATURES SHOULD EVER BE CHANGED IN THIS CLASS.

This class is used by GhidraScript.
Changing this class will break user scripts.
That is bad. Don't do that.



Program API



https://tenor.com/view/boom-mind-blown-mind-blowing-eyeglasses-gif-15569009

- Practically everything that is Ghidra
- All classes are available
- JavaDoc:

https://ghidra.re/ghidra_docs/api/index.html

• Source

https://github.com/NationalSecurityAgency/ghidra





Example 1: Listing imported functions

```
10 //This script lists all imported functions used by the binary.
  80 import ghidra.app.script.GhidraScript;
 10
    public class brucon 01 list imports extends GhidraScript {
 11
12
 130
        @Override
-14
        protected void run() throws Exception {
             String format = "%-30s %-30s\n":
15
             SymbolTable symbolTable = currentProgram.getSymbolTable();
16
17
             SymbolIterator symbolIter = symbolTable.getExternalSymbols();
 18
             printf("[-] Imported functions of %s\n\n", currentProgram.getExecutablePath());
             printf(format, "Function Name", "Library");
19
             printf("%s\n", "-".repeat(60));
 20
             for(Symbol symbol: symbolIter) {
 21
22
                 printf(format, symbol.getName(), symbol.getParentSymbol());
 23
 24
25
26 }
```



P-Code

- Intermediate Representation: lies between the assembly code and the decompiled code that the Ghidra UI shows
- Register transfer language, it translates every individual processor instruction to a sequence of P-Code operations
- For each instruction it describes the processor instruction, including all side effects



Raw P-Code

- This is what you see when you start printing P-Code from a Ghidra Script (we will see later)
- There are no higher-level connections (arguments, variables, etc...)

				CALLIND (FAM. 0X402048. 4)
0040100b	85 c0	TEST	EAX, EAX	
				(register, 0x200, 1) = COPY (const, 0x0, 1)
				(register, 0x20b, 1) = COPY (const, 0x0, 1)
				(unique, 0x42580, 4) = INT_AND (register, 0x0, 4), (reg
				(register, 0x207, 1) = INT_SLESS (unique, 0x42580, 4),
				(register, 0x206, 1) = INT_EQUAL (unique, 0x42580, 4),
				(unique, 0xd900, 4) = INT_AND (unique, 0x42580, 4), (co
				(unique, 0xd980, 1) = POPCOUNT (unique, 0xd900, 4)
				(unique, OxdaOO, 1) = INT_AND (unique, Oxd980, 1), (con
				(register, 0x202, 1) = INT_EQUAL (unique, 0xda00, 1), (
0040100d	76 76	11	LAB 004010	005

High P-Code

- Result of HighFunction()
- "High-level abstraction associated with a low-level function made up of assembly instructions. Based on information the decompiler has produced after working on a function."
- Closer to the C Pseudo Code
- Higher-level connections, such as function arguments, variables, etc...







Why use P-Code?

- Because it's cool
- Most of the binary analysis systems work with some kind of intermediate language.
- Compilers also use IR to separate the processes of parsing source code to a standardized format and translating the code to the machine code of the target architecture.
- Architecture independent
- Provides more information, then assembly (i.e. side effects)



https://pics.awwmemes.com/do-you-ever-think-happen-com-do-you-ever-think-screw-49053453.png





You get a Demo!

You get a Demo!



Everybody gets a DEMOOOOO!

Unrestricted | © Siemens 2021 | Gergely Revay | T RDA CST SES | 2021-10-07



Example 2: Print high P-Code at the COM functions' call site

- Malware often use Component Object Model (COM) as obfuscation
- COM offers standardized interfaces to various functionality
 - Control IE through COM
 - Control the firewall through COM
- Script identifies if binary uses COM objects
- Prints call sites for COM functions



Example 3: Recover CLSID and IID to identify which COM object is used

- CLSID: globally unique identifier of a COM class objects
- CLSIDs and their associated programs are recorded in the registry: HKEY_LOCAL_MACHINE\SOFTWARE\Classes\CLSID\{CLSID}
- To find out what COM object is used we need the CLSID
- IID: the interface that will be used to talk to the object.

\mathbf{C}	+	4	
- New Y			

HRESULT CoCreateInstance(
REFCLSID	rclsid,					
LPUNKNOWN	pUnkOuter,					
DWORD	dwClsContext,					
REFIID	riid,					
LPVOID	*ppv					
);						



Recap

- Ghidra is great
- Ghidra scripts and P-Code have a lot of potential for automation
- Ghidra is relatively well documented
- There are still challenges when one starts to dig deep
- Choose your tool for the task



https://www.avira.com/de/blog/nsa-macht-ghidra-oeffentlich-zugaenglich



References and Thanks

- Alexei Bulazel and Jeremy Blackthorne: <u>https://www.riverloopsecurity.com/blog/2019/05/pcode/</u>
- Rolf Rolles:
 <u>https://www.msreverseengineering.com/blog/2019/4/17/an-abstract-interpretation-based-deobfuscation-plugin-for-ghidra</u>
- Carlos Garcia Prado (@m0n0sapiens) helped with his infinite wisdom





Gergely Revay

LinkedIn: <u>https://www.linkedin.com/in/gergelyrevay/</u> Twitter: @geri_revay Youtube: <u>https://youtube.com/aetherlabnet</u>



